

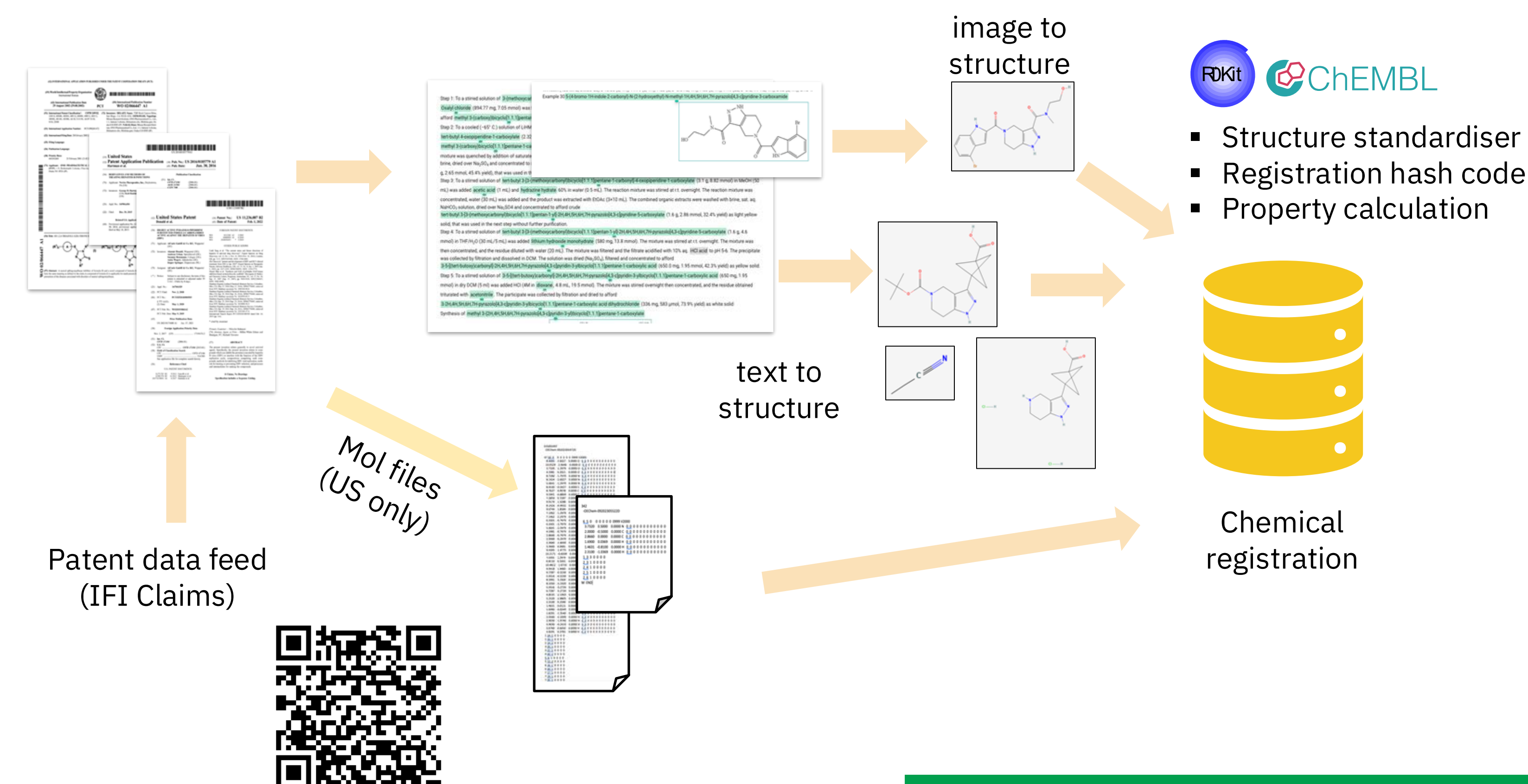
Chemical Biology Services, European Molecular Biology Laboratory-European Bioinformatics Institute (EMBL-EBI), Hinxton, Cambridge CB10 1SD, UK

The infographic displays four categories of data and a new system:

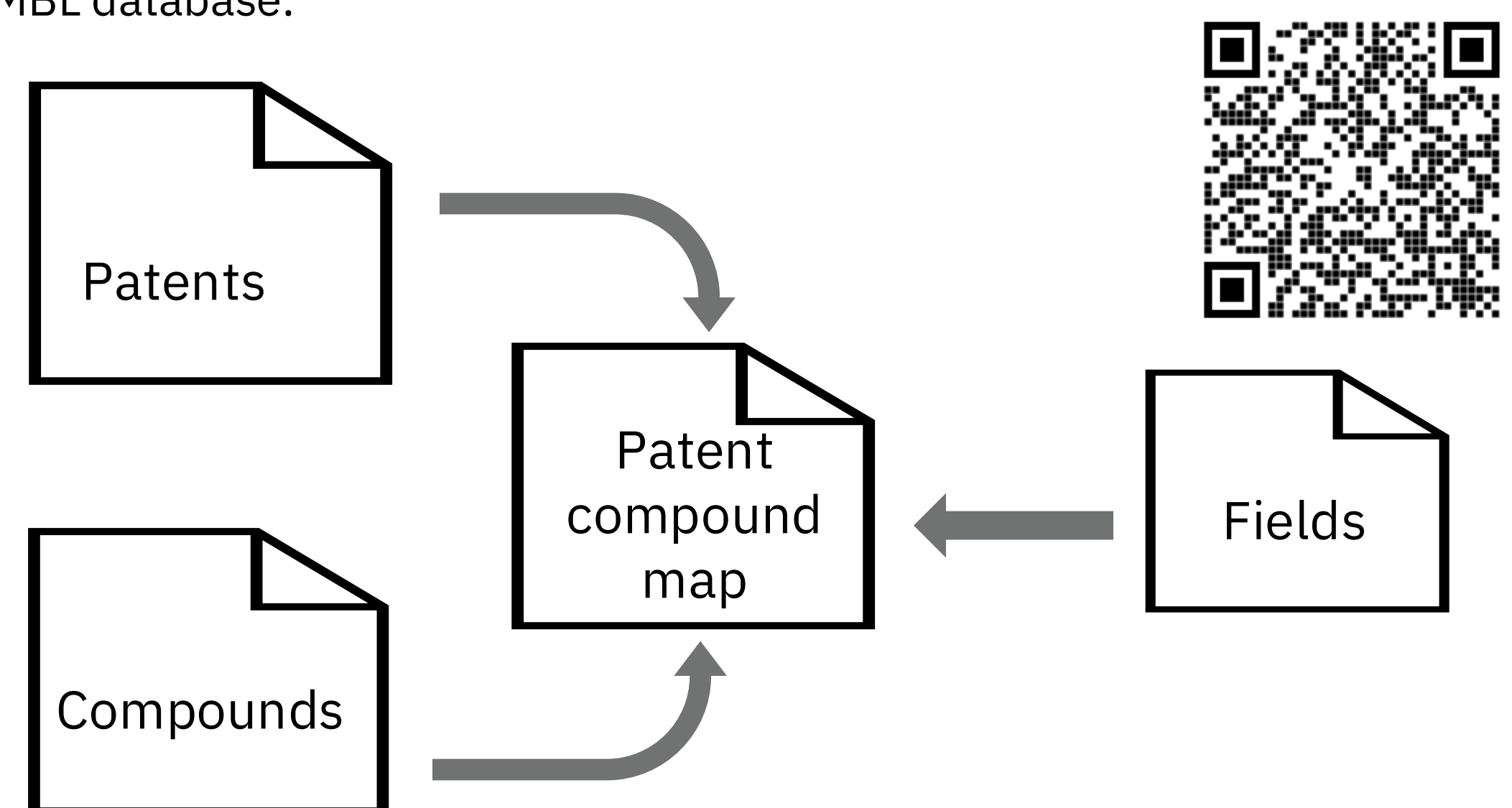
- 118M Patents**: Represented by a database cylinder icon.
- 43M Chemically-annotated documents**: Represented by a stack of documents icon.
- 31M Chemicals**: Represented by three chemical structures (benzene rings) icon.
- Biomedical annotation generated on the fly**: Represented by icons of a heart, a DNA helix, a virus, and a person.

A QR code is located in the bottom left corner, and a green banner at the bottom reads "New export system".

New export system




- Bulk data: new collection of core files representing the entire annotated SureChEMBL database.




- Data updates every two weeks
- Using highly compressed but still readable parquet files (direct query or insert into internal database)

www.surechembl.org



[Search](#)
[Downloads](#)
[Wiki](#)
[Contact Us](#)





Latest announcements

- Has **SureChEMBL** saved you time, effort or money?
- Please take 15 minutes to fill in our [impact survey](#) and help EMBL-EBI make the case for why open data resources are critical to life science research.
- Download the [SureChEMBL 2.0 data release update!](#)
- [SureChEMBL 2.0 announcement](#)

☒ All chemically annotated authorities
 ☐ Patents with small molecules
 ☐ Specify dates

SEARCH

 Query Assistant
  Structure Search

Total Hits: 37529

<

1

2

3

4

5

6

7

...

2,502

>

↓

Query: pa:bayer OR astra novartis OR genentech OR merck) AND desc:(chemotherap* AND ("Phosphoinoside 3 kinases"-3 OR PI3K)) AND (spectry:(US OR EP OR WO OR JP OR CN))

USE OF SUBSTITUTED 2,3-DIHYDROIMIDAZO[1,2-C]QUINAZOLINES

EP-2046049-A1

Substituted 2,3-dihydroimidazo[1,2-c]quinazoline derivatives useful for treating hyper-proliferative disorders and disease...

US-RE46556-E1


USE OF SUBSTITUTED 2,3-DIHYDROIMIDAZO[1,2-C]QUINAZOLINES FOR TREATING LYMPHOMAS

US-20170605873-A1

Main page with search by keyword returning a paginated patent list matching the query

[illegible]

Compound search result page



[Search](#)
[Downloads](#)
[Wiki](#)
[Contact Us](#)

Compound Details

SCEMBL3827

N-[(4-methyl-3-[[4-(pyridin-3-yl)pyrimidin-2-yl]amino]phenyl)-4-[(4-methylpiperidin-

SMILES: CN1CCN(CC(C)=O)C=C(C)C=CN2C=CN(C)C(NC3=CC=CC=C3OC(=O)C=C2)C1

Amble: HCN=H[C@@H](N)O

C1=21-15251-18272138-29-31-13112423-2924-4-5-12-30-19-24322823723-8-6-227-9-232036-16-14-352315-17-36-43-13-18-1961-16-1720402-329313237701312834

InChI Key: UTPFNSKRBQBMGHWUHFFFAZYSN.N

Log P: 4.59

Mol Weight: 493.60

UniChem Cross References	▼
Patents for compound	▲
Total patents found: 98694	<div> < 1 2 3 4 5 6 7 8 9 ... 9,870 > </div>
Anti-cancer sustained-released implantation agent	
CN-10138399B-A	
Amide derivatives and their application for the treatment of g protein related diseases	
CN-101405268-A	
Effervescent tablet containing imatinib mesylate and preparation method thereof	
CN-1014601797-A	
Sustained-release preparation containing chemotherapy synergist for treating solid tumors	
CN-101444483-A	

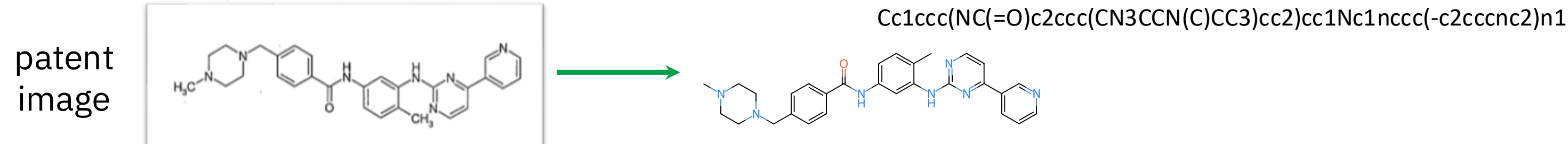
Compound details page with compound structure information, UniChem cross referencing, and list of patents where the compound was found

- Convert image to molecule structure
- molecule structure once converted

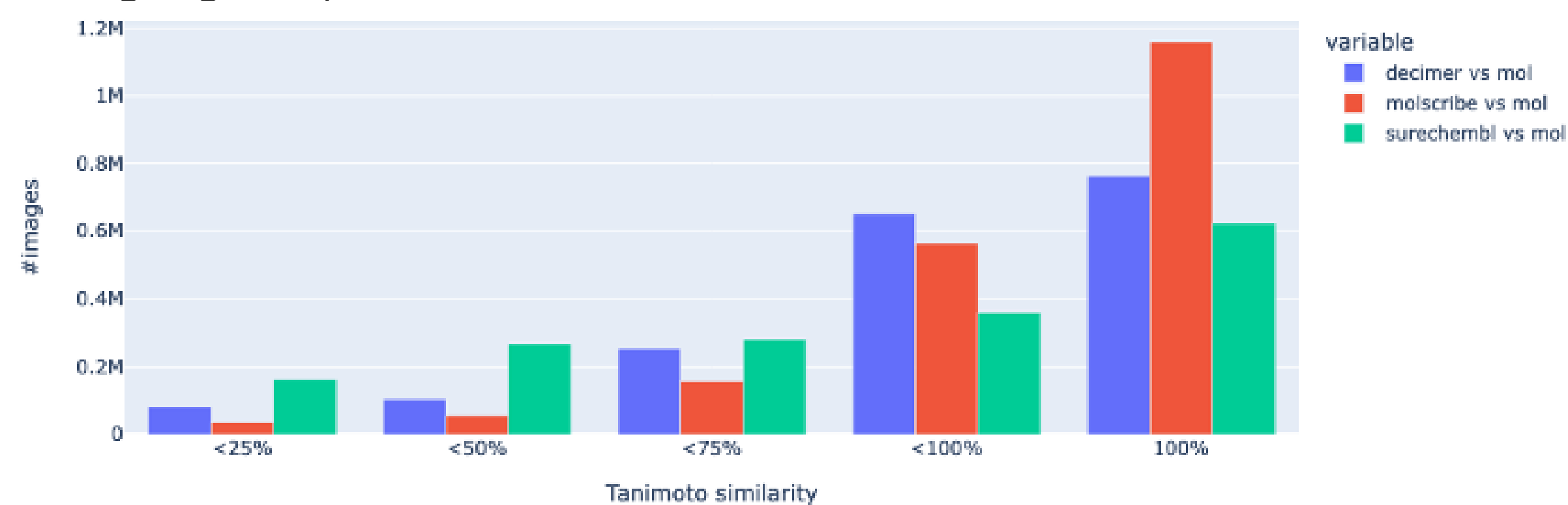
```

Cc1ccc(NC(=O)c2ccc(CN3CCN(C)CC3)cc2)cc1Nc1nccc(-c2ccccc2)n1

```



- New generation OCR tools were evaluated on 2023 US patent dataset: DECIMER [4] / MolScribe [5]
- Compare structures from MOL files with structures extracted from associated images
- AI-based model, in particular MolScribe, gives improved structure accuracy from image extraction
- Now on-going incorporation of MolScribe into SureChEMBL2.0

[illegible][illegible]

Patent annotation with Leadnine (prototype), orange: generic chemical name, pink: generic molecule, grey: anatomy, violet: molecule dictionary, turquoise: mechanism, green: PubChem dictionary, dark red: gene, yellow: polymer, light red: journal, khaki: organism, dark orange: disease

[1] Senger, S. (2017). Assessment of the significance of patent-derived information for the early identification of compound–target interaction hypotheses. *J. Cheminform.* 9, 26.

[2] Papadatos, G., Davies, M., Dedman, N., Chambers, J., Gaulton, A., Siddle, J., Koks, R., Irvine, S.A., Pettersson, J., Goncharoff, N., et al. (2016). SureChEMBL: a large-scale, chemically annotated patent document database. *Nucleic Acids Res.* 44, D1220–D1228

[3] Lowe, D., Sayle, R. (2015). LeadMine: a grammar and dictionary driven approach to entity recognition. J. Cheminform. 7, S1, S5

[4] Rajan, K., Brinkhaus, H.O., Agea, M.I., Zieslesny, A., and Steinbeck, C. (2023). DECIMER.ai: an open platform for automated optical chemical structure identification, segmentation and recognition in scientific publications. *Nat. Commun.* **14**, 5045.

[5] Qian, Y., Guo, J., Tu, Z., Li, Z., Coley, C.W., and Barzilay, R. (2023). MolScribe: Robust Molecular Structure Recognition with Image-to-Graph Generation. *J. Chem. Inf. Model.* 63, 1925–1934