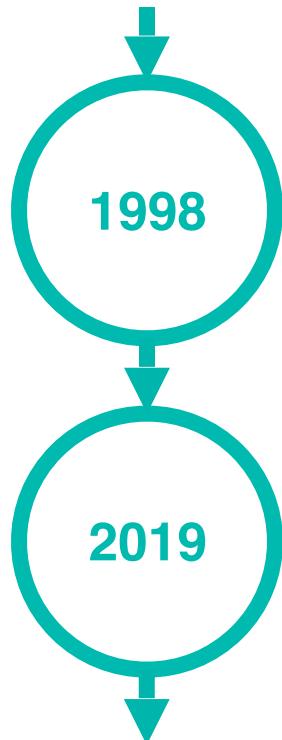


CHEMICAL SEARCH WITH NEO4J AND KNOWLEDGE EXTRACTION FROM SCIENTIFIC PAPERS

ChemAxon





HQ - Budapest (Hungary)



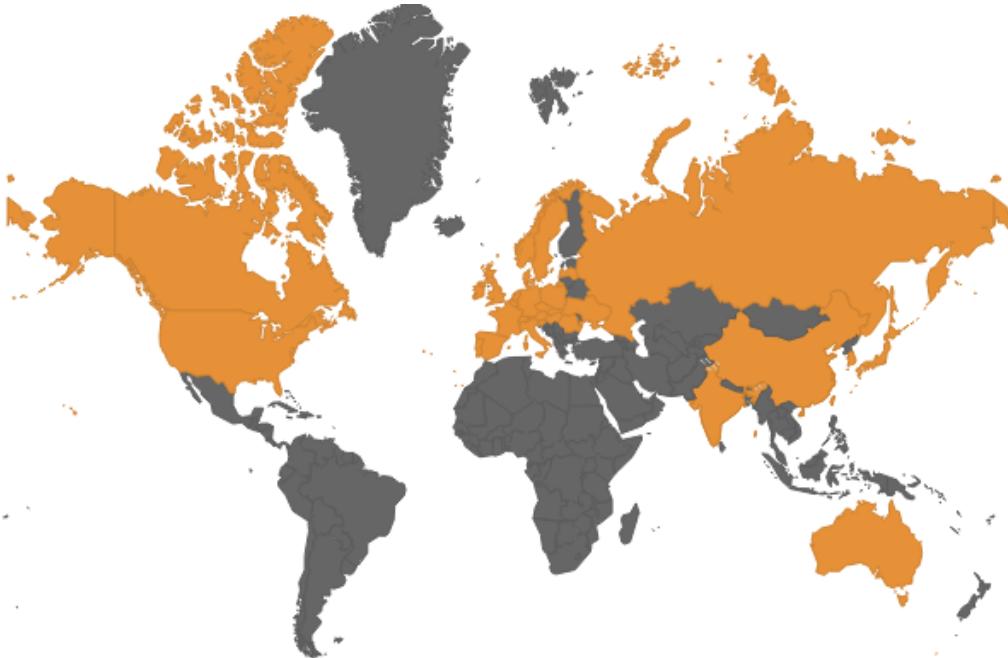
**Development,
Consultancy**

Prague (Czech Republic)



**Sales, Consultancy and Application Scientists
Boston and San Diego (United States)**

ChemAxon around the Globe



Distributors Spain, India, China, Korea, Japan, Singapore, Australia

Clients 500+

Partners 50+

Users 300 000+

The background features a complex, abstract geometric pattern composed of numerous thin, semi-transparent purple lines and small purple dots. These lines form a dense network of triangles and other polygons, creating a sense of depth and connectivity. The pattern is concentrated in the lower-left quadrant of the slide, with some lines extending towards the center.

THE PROJECT

Input data set

The screenshot shows the homepage of the US National Library of Medicine's PubMed Central (PMC) site. At the top, there is a navigation bar with links for NCBI, Resources, How To, and Sign in to NCBI. Below the navigation bar is the PMC logo and the text "US National Library of Medicine, National Institutes of Health". A search bar is prominently displayed, with "PMC" selected from a dropdown menu. To the right of the search bar is a "Search" button. Below the search bar are links for "Advanced" and "Journal list". At the bottom of the header are links for "About PMC", "For Publishers", and "Related Resources".

Open Access Subset

The PMC Open Access Subset [OA](#) is a part of the total collection of articles in PMC. The articles in the OA Subset are made available under a Creative Commons or similar license that generally allows more liberal redistribution and reuse than a traditional copyrighted work.

To preview the articles or get a current count of articles in the OA Subset, do a search for [open access\[filter\]](#) in PMC. As of May 2018, there were over 2 million articles available in this collection.

<ftp://ftp.ncbi.nlm.nih.gov/pub/pmc/>

Data set

INPUT

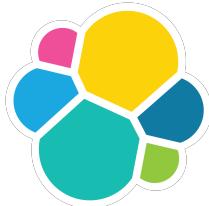
- Documents from <ftp://ftp.ncbi.nlm.nih.gov/pub/pmc/>
- size: **~30 Gb** (compressed)
- format: **xml**
- Document count: **~1.9 million**

OUTPUT

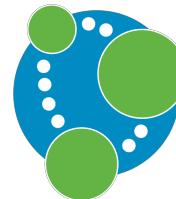
- Chemical database
 - Distinct molecules: **~211.000**
 - All occurrences: **~53 million**
 - MSSQL DB: 800 GB
- Graph DB: 10.5 GB
- Free text DB: 110 GB

Application Status	
Working extractors	12
Database status	✓
Content sources	5
Overall document count	1876605
Distinct structure count	211706
All molecule hits	53488909

Tools



elastic



neo4j

HARDWARE

PROVIDER	TYPE	MEMORY	CPU COUNT	DISK SIZE
Amazon EC2	M4.4xlarge	64 GB	16	1.5 TB

CHEMICAL DATA EXTRACTION

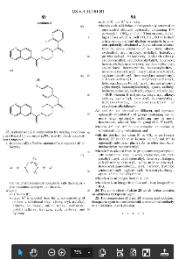
Chemical files and formats

OPTICAL CHARACTER RECOGNITION		MRV	PPTX		
RXN	SDF	OPTICAL STRUCTURE RECOGNITION			
CDXML	IUPAC	CAS	XLS	CDX	SKC
MARVIN OLE	ACCORD FOR EXCEL		DOC	MOL	
JCHEM FOR EXCEL	DOCK	PPT	ISIS OLE		
XLSX	JCHEM FOR SHAREPOINT		ISIS FOR EXCEL		
CORPORATE ID		SMILES	CHEMDRAW OLE		
RDF	CML	INCHI	SMARTS	CXSMILES	
CXSMARTS		XYZ	COMMON NAME		ONE
SYSTEMATIC NAME			CHEMDRAW FOR EXCEL		

Corporate ID

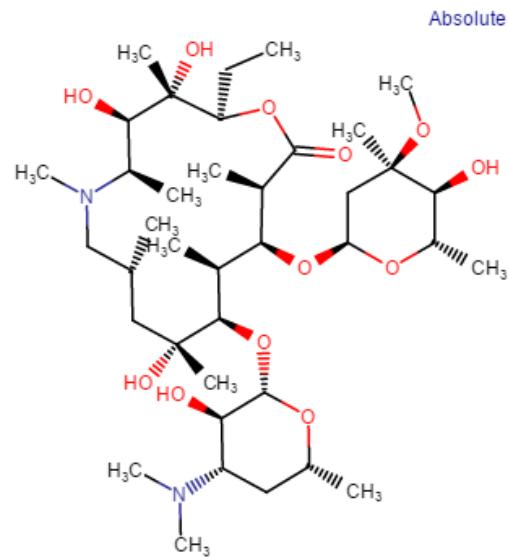
SMDL-00000420
SMDL-00000421
SMDL-00000422
SMDL-00000423
SMDL-00000424
...
SMDL-00116392
SMDL-00116393

OSR / OCR

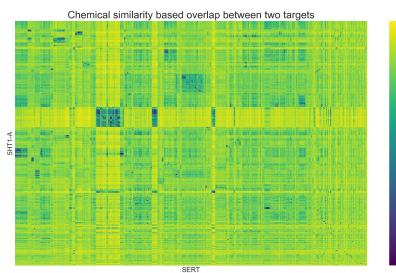
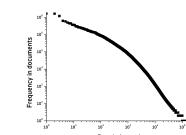
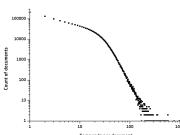
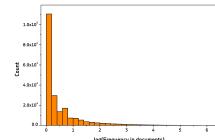
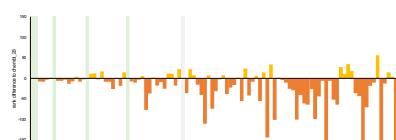
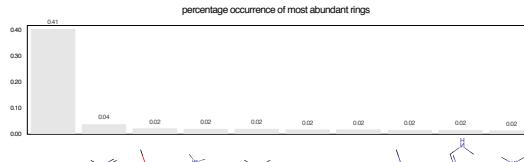
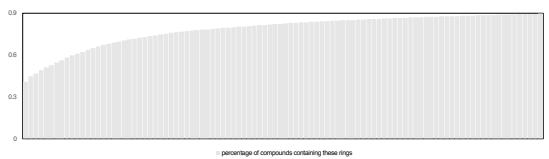
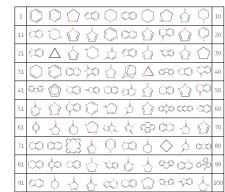
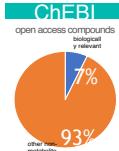
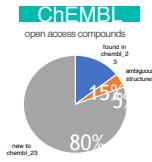
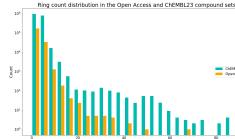
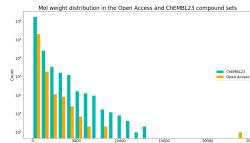
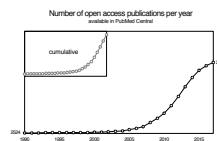
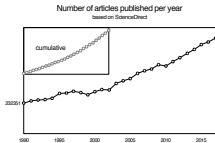


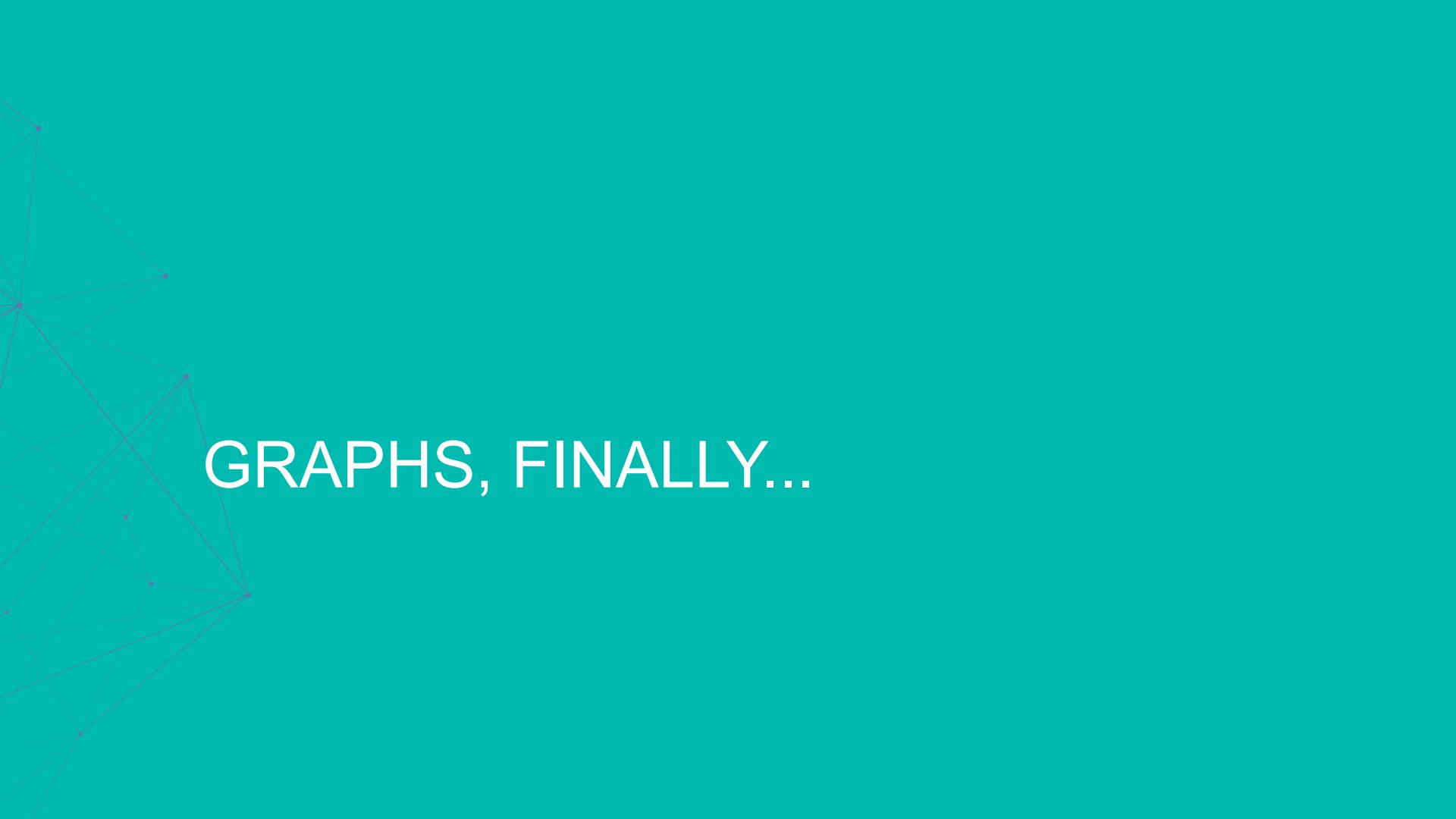
Nomenclature

- Azadose
Azithromycin
Azithromycin Dihydrate
Azithromycin Monohydrate
Azithromycin Pfizer Brand
Aztreonam
Azythromycin
Bayer Brand of Azithromycin Dihydrate
Dihydrate, Azithromycin
Funk Brand of Azithromycin Dihydrate
Goxal
Lesvi Brand of Azithromycin Dihydrate
Mack Brand of Azithromycin Dihydrate
Monohydrate, Azithromycin
Pfizer Brand of Azithromycin
Pfizer Brand of Azithromycin Dihydrate
Pharmacia Brand of Azithromycin Dihydrate
Sumamed
Toraseptol
Ultreon
Vinzam
Vita Brand of Azithromycin Dihydrate
Zenaventon
Zithromax
Zitromax



CONTENT ANALYSIS



The background features a complex, abstract geometric pattern composed of numerous thin, semi-transparent purple lines and small purple dots. These lines form a dense network of triangles and other polygons across the entire slide.

GRAPHS, FINALLY...

Graph path finding

DB size

10.42 GB

Nodes count

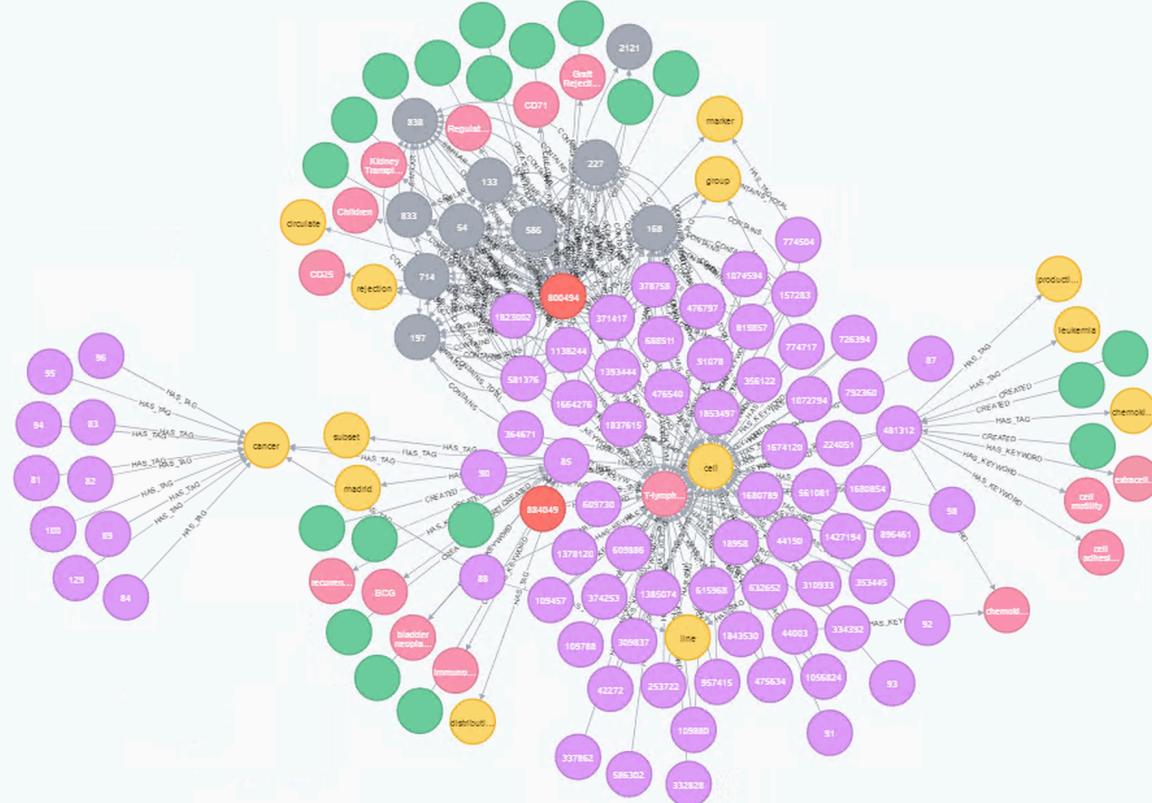
13 769 389

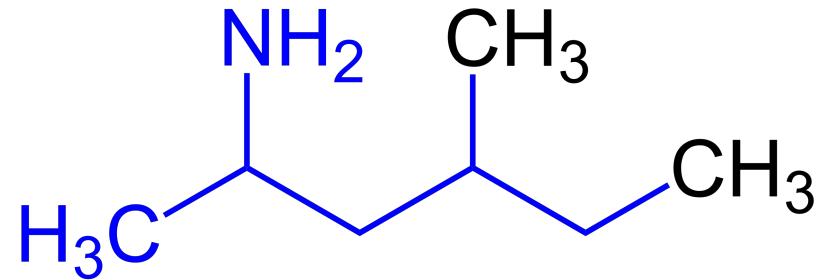
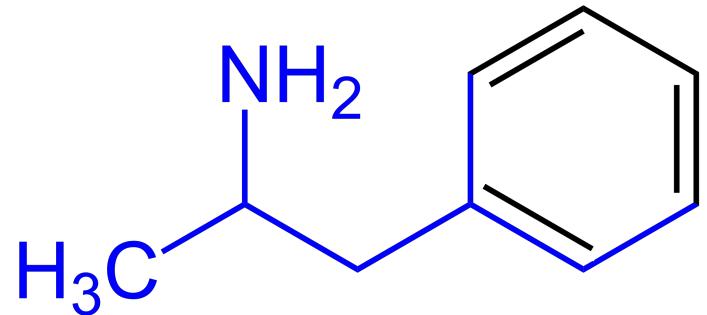
Relationships

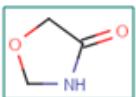
124 442 817

Relationship types

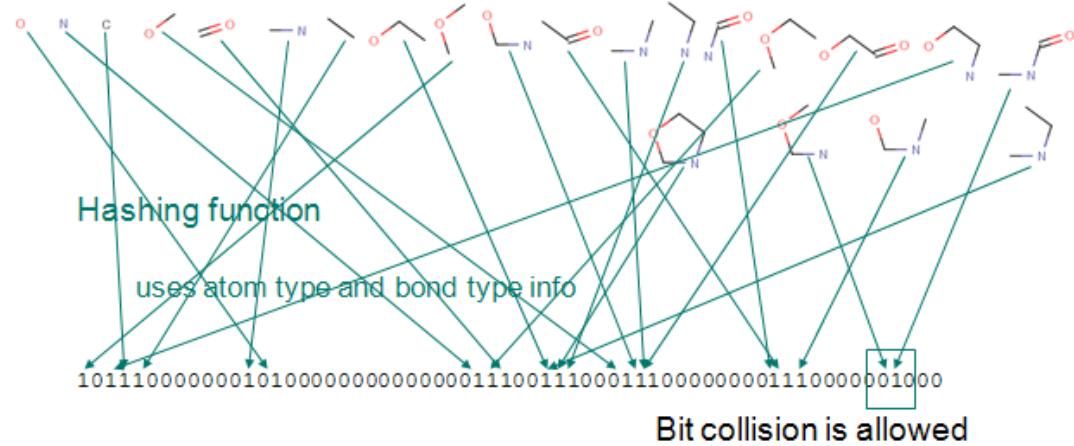
**Keys, tags,
similarity,
contains,
references,
created (by
authors)**







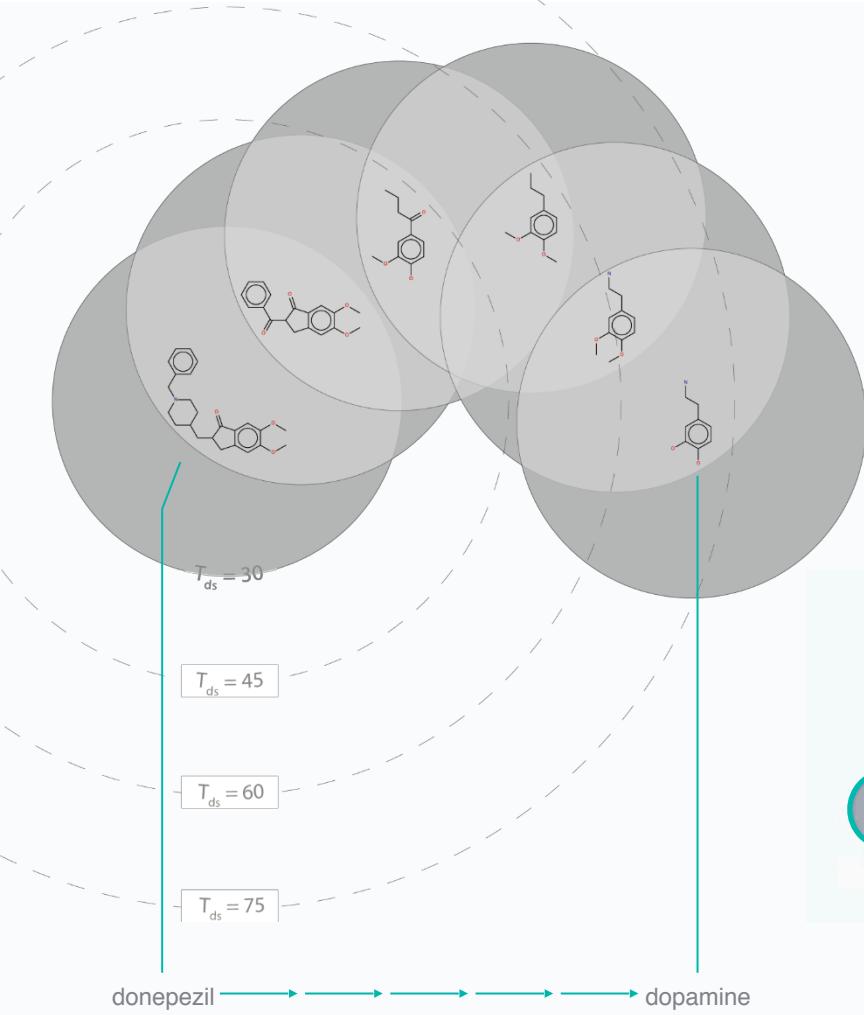
Patterns in the molecule (Note – all substructures!):



Chemical hashed fingerprint

$$T_s(X, Y) = \frac{\sum_i (X_i \wedge Y_i)}{\sum_i (X_i \vee Y_i)}$$

Tanimoto similarity



GOAL

- Find a route from one chemical space to another hopping through molecules
- In maximum 5 steps

~8 sec

