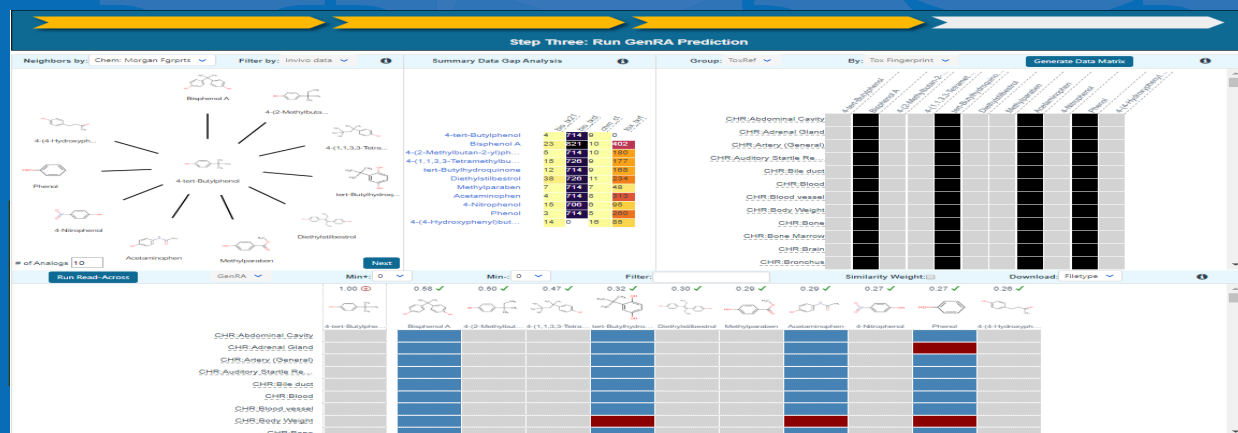


Building Scientific Confidence in the Development and Application of Objective Read-across Approaches



Grace Patlewicz
Center for Computational Toxicology & Exposure (CCTE), US EPA

Acknowledgements

- Imran Shah - co-lead on Generalised Read-across (GenRA)
- George Helman (former student)
- Tia Tate
- Willysha Jenkins

Outline

- Read-across - definition
- Background context - tools, frameworks
- Generalised Read-across (GenRA): A data driven approach to read-across
 - Implementations of GenRA
 - Recent work applying GenRA
- ICCVAM Read-Across Workgroup activities

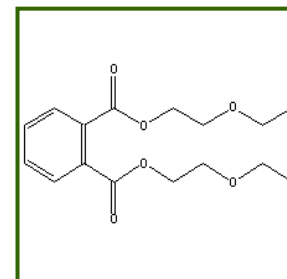
Read-across

- Read-across describes the method of filling a data gap whereby a chemical with existing data values is used to make a prediction for a 'similar' chemical.
- A target chemical is a chemical which has a data gap that needs to be filled i.e. the subject of the read-across.
- A source analogue is a chemical that has been identified as an appropriate chemical for use in a read-across based on similarity to the target chemical and existence of relevant data.

	Source chemical	Target chemical
Property	●	○

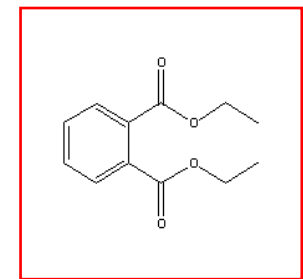
● Reliable data

○ Missing data



Known to be
harmful

Acute
toxicity?



Predicted to be
harmful

Ongoing issues with read-across

- Although there is much guidance for developing read-across assessment, acceptance remains an issue, not helped since read-across still remains a subjective, expert driven assessment.
- One issue thwarting acceptance relates to the “uncertainty of the read-across prediction”.
- As such there have been many efforts to identify the sources of uncertainty in read-across, characterise them in a consistent manner and identify practical strategies to address and reduce those uncertainties.
- Notable in these efforts have been the development of frameworks for the assessment of read-across, evaluating the utility of New Approach Methods (NAMs).
- Quantifying uncertainty and performance of read-across is still a need as are ways to better characterise different similarity contexts (metabolism, reactivity etc.)

Read-Across Tools



Navigating through the minefield of read-across tools: A review of in silico tools for grouping

Grace Patlewicz^{a,*}, George Helman^{a,b}, Prachi Pradeep^{a,b}, Imran Shah^a

^a National Center for Computational Toxicology (NCCT), Office of Research and Development, US Environmental Protection Agency, 109 TW Alexander Dr, Research Triangle Park (RTP), NC 27711, USA

^b Oak Ridge Institute for Science and Education (ORISE), Oak Ridge, TN, USA

ARTICLE INFO

Article history:
Received 29 March 2017
Received in revised form 22 May 2017
Accepted 25 May 2017
Available online 29 May 2017

Keywords:
Category approach
Analogue approach
Data gap filling
Read-across
(Q)SAR
Trend analysis
Nearest neighbours

ABSTRACT

Read-across is a popular data gap filling technique used within analogue and category a regulatory purposes. In recent years there have been many efforts focused on the challenge in read-across development, its scientific justification and documentation. Tools have also been developed to facilitate read-across development and application. Here, we describe a number of available read-across tools in the context of the category/analogue workflow and review their capabilities, strengths and weaknesses. No single tool addresses all aspects of the workflow, how the different tools complement each other and some of the opportunities for their future development to address the continued evolution of read-across.

Published by



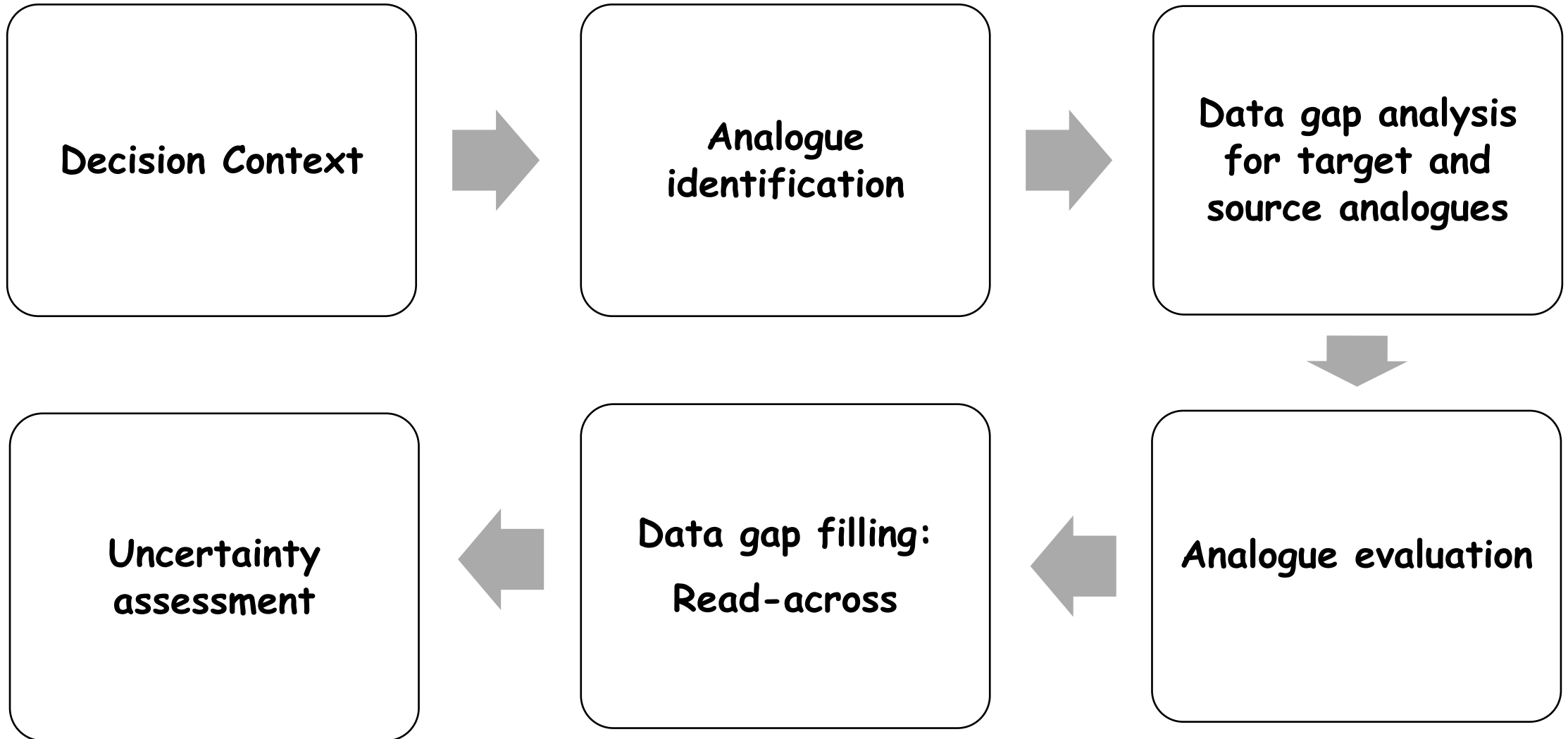
Summary of key features of selected publicly available read-across tools.

	AIM	ToxMatch	Ambit	OECD Toolbox	CBRA	ToxRead	CIPro
Development timeline	Java based version is dated 2012. Initial development of web version was 2005.	First public version released in Dec 2006	Original AMBIT tool was developed in 2004–2005	Proof of concept released in 2008	Implementation of the Low et al. [27]	Implementation of Gini et al. [22]	Implementation described in Russo et al. [45]
Type of Tool	Standalone	Standalone	Web-based and standalone	Standalone or Client/Server	Standalone	Standalone	Web-based
Latest Version	1.01 (Nov 2013) Static	1.07 (Jan 2009) Static	3.0.3 Ongoing Enhanced in 2013–2015	3.4 (July 2016) Version 4 released April 2017 Ongoing	0.75 First release	0.11 BETA Ongoing	First release
Developed by	SRC Inc	Ideaconsult Ltd	Ideaconsult Ltd	LMC, Bourgas	Fourches Lab at North Carolina State University	Istituto di Ricerche Farmacologiche Mario Negri http://www.toxread.eu/	Zhu Research Group at Rutgers University http://ciipro.rutgers.edu/
Available from	https://www.epa.gov/tscascreening-tools/analogue-identification-methodology-aim-tool	https://eur1-ecvam-jrc.ec.europa.eu/laboratories-research/predictive-toxicology/qsar_tools/toxmatch	http://cetic-lri.org/lri_toolbox/ambit/	www.qsartoolbox.org	http://www.fourches-laboratory.com/software	http://www.toxread.eu/	http://ciipro.rutgers.edu/
Accepted Chemical Input	CAS, Name, SMILES, structure drawing/import	CAS, Name, SMILES, InChI	Name, identifiers, SMILES, InChI	CAS, Name, SMILES, structure drawing, MOL, sdf	Mol file, descriptors as txt	SMILES	PubChem CID, CAS, IUPAC, SMILES, InChI
Endpoint Coverage	N/A	Any based on user input	IUCLID ^a 5-supported endpoints (43 total)	Any as per the regulatory endpoints	Any based on user input	Mutagenicity and Bioconcentration Factor (BCF)	Any based on user input
Analogue Identification Approach	Fragment matching	Distance and correlation based similarity indices based on descriptors or fingerprints	Substructure or similarity searching using structure, name, SMILES, InChI Manual	Category definition followed by subcategorisations	Tanimoto distance using chemical and biological descriptors	VEGA similarity algorithm	Weighted Estimated Biological Similarity
Neighbour Selection	Automatic	Automatic	Automatic	Automatic + Manual Filter	Automatic	Automatic	Automatic + Manual Filter
Data Source	Tool provides inventory index	User provided or tool provided	User and tool provided	User provided or tool provided	User provided	Tool provided as a result of the EU ANTARES project	User provided but tool provides PubChem in vitro data
Quantitative vs Qualitative	N/A	Both	User determined – Qualitative	Both	Qualitative	Qualitative for mutagenicity, quantitative for BCF	Qualitative
Visualisation	None	Standard 2D plots, histograms and similarity matrix	None	Standard 2D Plots	Radial plot of neighbours	Interactive Neighbour plot	Activity Plot
Output/Export	Output reports in the form of HTML, pdf or Excel	sdf or txt files of data, image files of plots	Assessment report as docx or xlsx, data matrix as xlsx	IUCLID format, pdf and rtf files of prediction report, text files of data, image files of plots etc	NA	Image file of plot	Tabulation of predictions and image of similarity plot

^a IUCLID stands for International Uniform Chemical Information Database. IUCLID is a software program for the administration of data on chemical substances first developed to fulfill EU information requirements under REACH.

(Patlewicz et al., 2017)

Read-across workflow



A harmonised hybrid read-across workflow



Contents lists available at ScienceDirect

Computational Toxicology

journal homepage: www.elsevier.com

Journal
Cover
Image

Navigating through the minefield of read-across frameworks: A commentary perspective

Grace Patlewicz^{a,*}, Mark T.D. Cronin^b, George Helman^{a,c}, Jason C. Lambert^d, Lucina E. Lizarraga^d, Imran Shah^a

^a National Center for Computational Toxicology (NCCT), Office of Research and Development, US Environmental Protection Agency (US EPA), 109 TW Alexander Dr, Research Triangle Park (RTP), NC 27711, USA

^b School of Pharmacy and Biomolecular Sciences, Liverpool John Moores University, Byrom Street, Liverpool L3 3AF, UK

^c Oak Ridge Institute for Science and Education (ORISE), 1299 Bethel Valley Road, Oak Ridge, TN 37830, USA

^d National Center for Evaluation Assessment (NCEA), US Environmental Protection Agency (US EPA), 26 West Martin Luther King Dr, Cincinnati, OH 45268, USA

- Where do NAM data fit?
- How should we transition to data-driven approaches?
- Quantifying the uncertainty in the read-across predictions made?

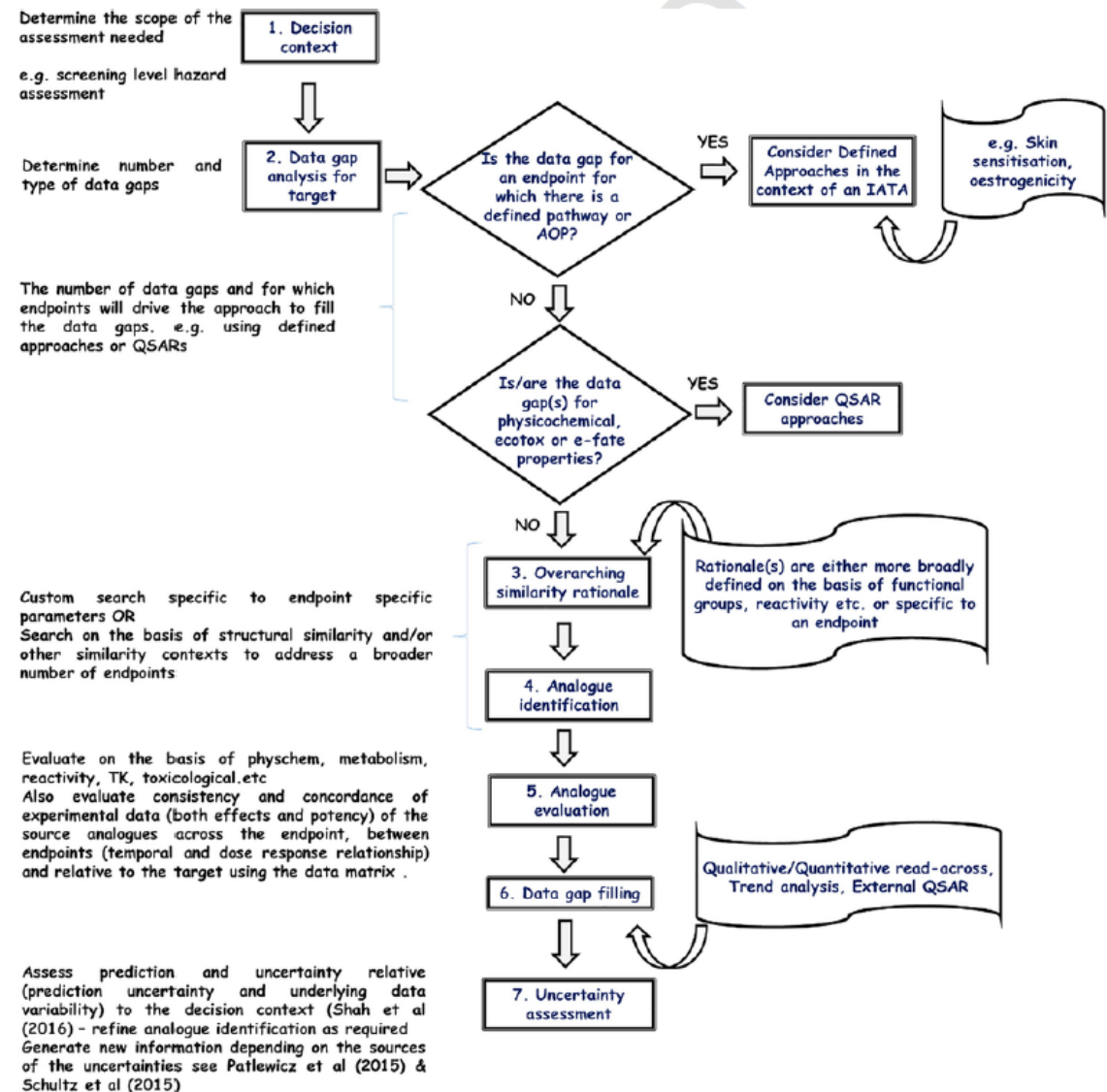


Fig. 9. A harmonised hybrid development and assessment framework.

GenRA (Generalised Read-Across)

- Predicting toxicity as a similarity-weighted activity of nearest neighbours based on chemistry and bioactivity descriptors (Shah et al, 2016)
- Goal: To establish an objective performance baseline for read-across and quantify the uncertainty in the predictions made

$$y_i^{\beta, \alpha} = \frac{\sum_j^k s_{ij}^{\alpha} x_j^{\beta}}{\sum_j^k s_{ij}^{\alpha}}$$

Jaccard similarity:

$$s_{ij} = \frac{\sum_l (x_{il} \wedge x_{jl})}{\sum_l (x_{il} \vee x_{jl})}$$

$\alpha \in \{chm, bio, bc\}$

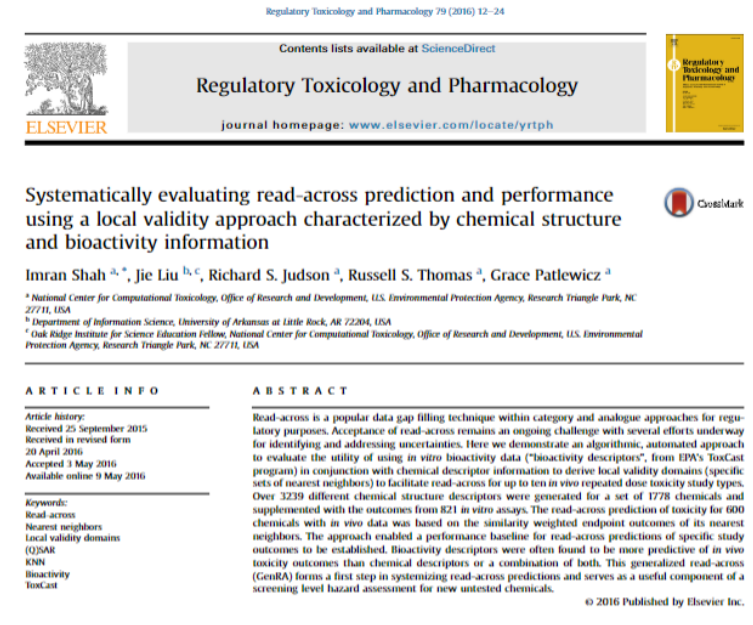
$\beta \in \{bio, tox\}$

$y_i = \text{predicted activity of chemical}(c_i)$

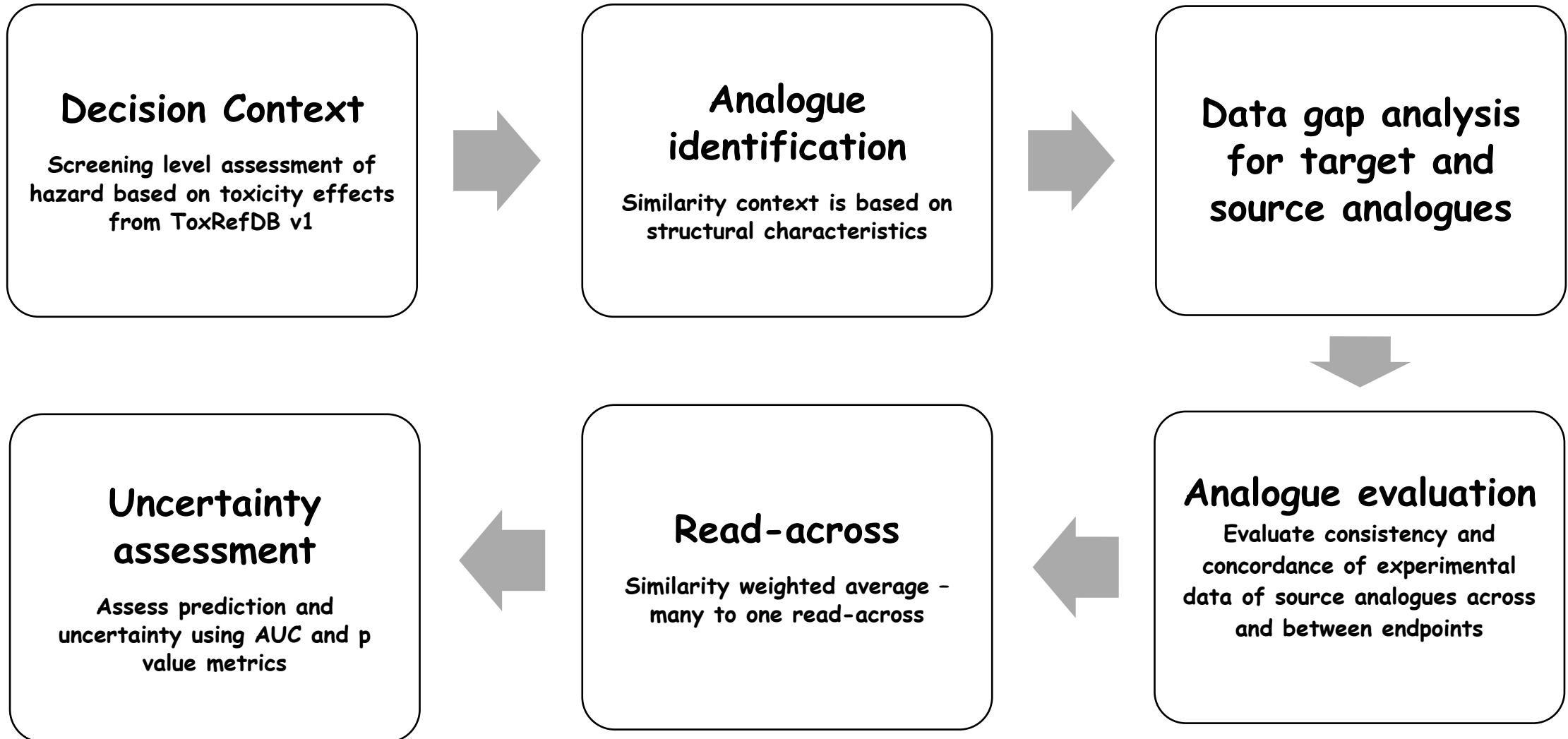
$x_j^{\beta} = \text{activity of } c_j \text{ in } \beta$

$s_{ij}^{\alpha} = \text{Jaccard similarity between } x_i^{\alpha}, x_j^{\alpha}$

$k = \text{up to } k \text{ nearest neighbours}$

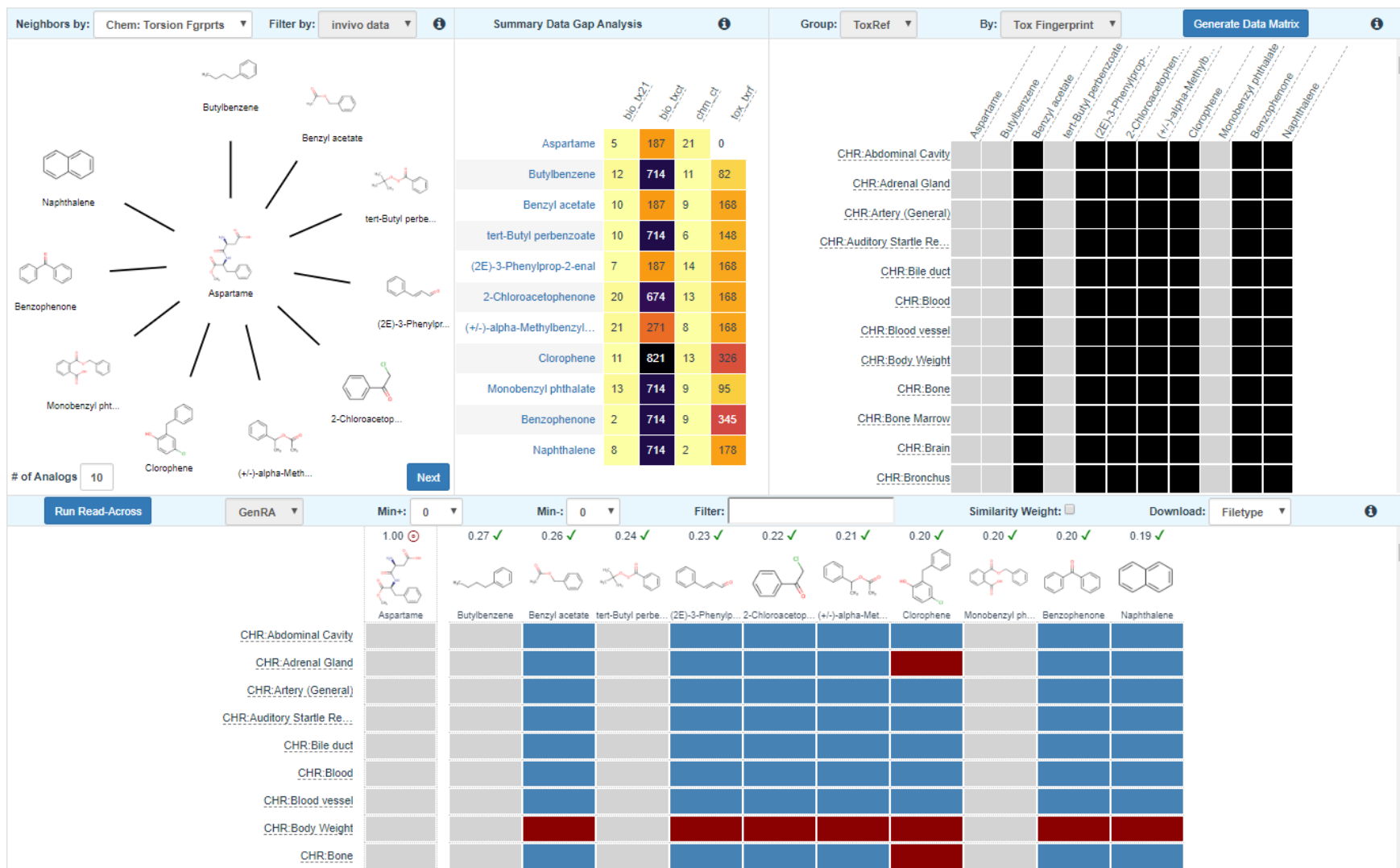


Read-across workflow in GenRA v1.0

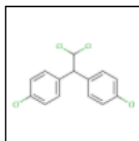


GenRA tool in reality

- GenRA v1.0 Integrated into the EPA CompTox Chemicals Dashboard



GenRA tool in practice



p,p'-DDD

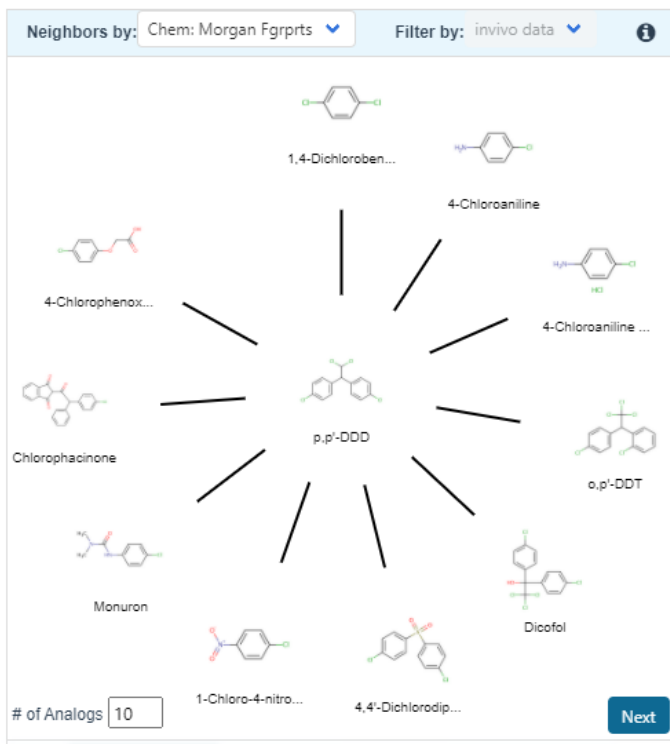
72-54-8 | DTXSID4020373

Searched by DSSTox Substance Id.

Search for a chemical and click on the GENRA link on the lefthand panel

Generalized Read-Across (GenRA)

Step One: Analog Identification and Evaluation

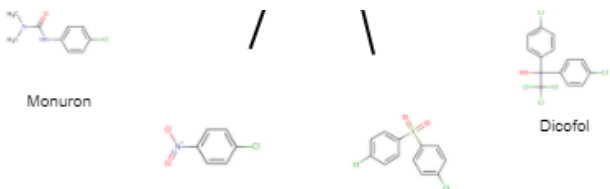


Radial plot with Target chemical of interest in the centre and source analogues (similar) ordered clockwise by decreasing similarity (Jaccard)

GenRA tool in practice

Step Two: Data Gap Analysis & Generate Data Matrix

- How data poor is my target and what data exists for the source analogues identified
- Do they address the data gaps of interest for the target chemical?



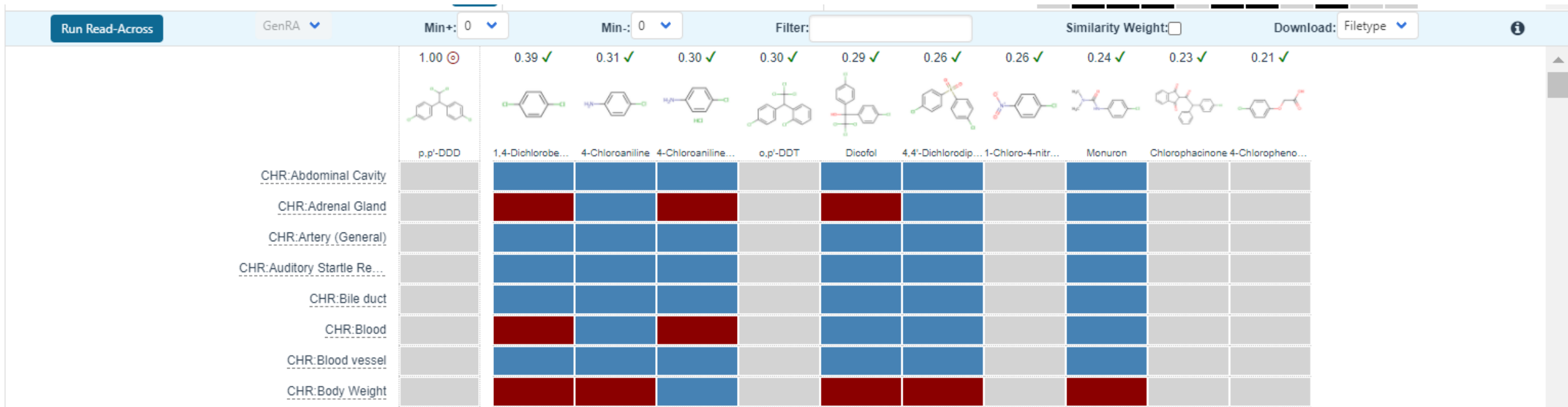
	bio tx21	bio txct	chm ct	tox txrf
p,p'-DDD	42	714	10	0
1,4-Dichlorobenzene	7	714	4	345
4-Chloroaniline	6	714	6	83
-Chloroaniline hydrochl...	17	0	7	168
o,p'-DDT	37	726	12	177
Dicofol	40	818	17	345
,4'-Dichlorodiphenyl sul...	9	271	5	168
1-Chloro-4-nitrobenzene	10	674	5	167
Monuron	12	714	7	168
Chlorophacinone	51	234	19	95
4-Chlorophenoxyacetic ...	9	232	8	180

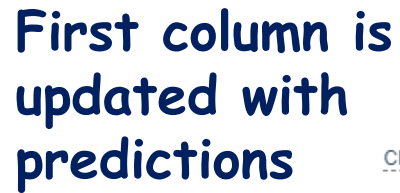
	p,p'-DDD	1,4-Dichlorobenzene	4-Chloroaniline	4-Chloroaniline hydr...	o,p'-DDT	Dicofol	4,4'-Dichlorodiphen...	1-Chloro-4-nitroben...	Monuron	Chlorophacinone	4-Chlorophenoxyac...
CHR:Abdominal Cavity											
CHR:Adrenal Gland											
CHR:Artery (General)											
CHR:Auditory Startle Re...											
CHR:Bile duct											
CHR:Blood											
CHR:Blood vessel											
CHR:Body Weight											
CHR:Bone											
CHR:Bone Marrow											
CHR:Brain											
CHR:Bronchus											

Next

Should I consider subcategorising the analogues selected?

Toxicity data represented as binary outcomes - red (positive), blue (negative), grey (no data)





GenRA Tool in practice

- Database underpinning GenRA v1.0: ToxRefDB v1
 - Different study types and effects within them are predicted e.g. chronic_liver is annotated as CHR_liver
 - Negative results - assume that if a particular guideline study was conducted but the effects were not reported than a chemical would be negative for that particular effect for that type of guideline study
 - Positive results - min dose at which toxicity effects are observed in a study
- Prediction: Similarity weighted activity
- Performance is categorised by the AUC of the ROC
 - The significance was empirically estimated by constructing a null distribution by permuting the toxicity values 100 times and calculating the fraction of times the AUC was more extreme than what would be observed by chance (this is reported as the p-value).

- [illegible]

GenRA Tool in practice

- Rank order positive results based on AUC and p values
- Look at the distribution of positive vs negatives predictions
- Explore what effects are being identified for the source analogues - consider identifying the underlying data for source analogues (elsewhere on the Dashboard) - is there a critical effect that is driving the toxicity that should be compared with the target chemical predictions?
-
- Depends on the decision context and the level of uncertainty that can be tolerated.

GenRA tools

- Efforts are underway to update the underlying data sources of the webapp GenRA for a summer release*
- An alternative and programmatic batch means of using GenRA is available through `genra-py`*, a standalone python library to enable user specific datasets to be analysed - see <https://github.com/i-shah/genra-py> (Shah et al, 2021)

Bioinformatics, 2021, 1–2
doi: 10.1093/bioinformatics/btab210
Advance Access Publication Date: 27 March 2021
Applications Note

OXFORD

Data and text mining

Generalized Read-Across prediction using `genra-py`

Imran Shah , Tia Tate and Grace Patlewicz

Center for Computational Toxicology and Exposure, Office of Research and Development, United States Environmental Protection Agency, Research Triangle Park, NC 27709, USA

*To whom correspondence should be addressed.
Associate Editor: Jonathan Wren

Received on December 14, 2020; revised on March 15, 2021; editorial decision on March 24, 2021; accepted on March 25, 2021

Abstract

Motivation: Generalized Read-Across (GenRA) is a data-driven approach to estimate physico-chemical, biological or eco-toxicological properties of chemicals by inference from analogues. GenRA attempts to mimic a human expert's manual read-across reasoning for filling data gaps about new chemicals from known chemicals with an interpretable and automated approach based on nearest-neighbors. A key objective of GenRA is to systematically explore different choices of input data selection and neighborhood definition to objectively evaluate predictive performance of automated read-across estimates of chemical properties.

Results: We have implemented `genra-py` as a python package that can be freely used for chemical safety analysis and risk assessment applications. Automated read-across prediction in `genra-py` conforms to the scikit-learn machine learning library's estimator design pattern, making it easy to use and integrate in computational pipelines. We demonstrate the data-driven application of `genra-py` to address two key human health risk assessment problems namely: hazard identification and point of departure estimation.

Availability and implementation: The package is available from github.com/i-shah/genra-py.

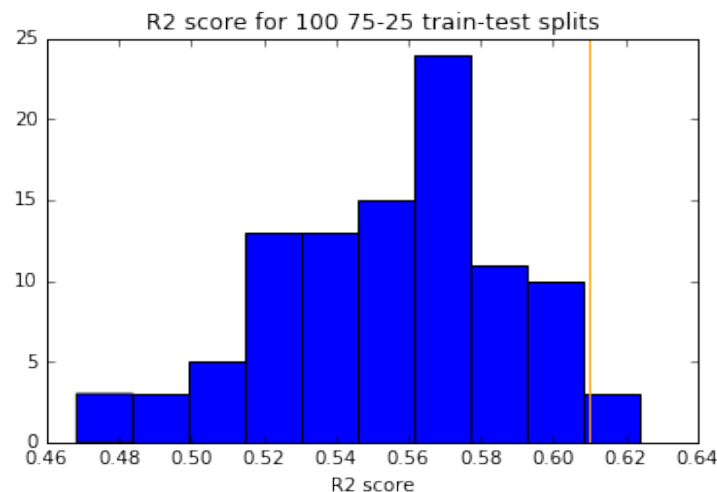
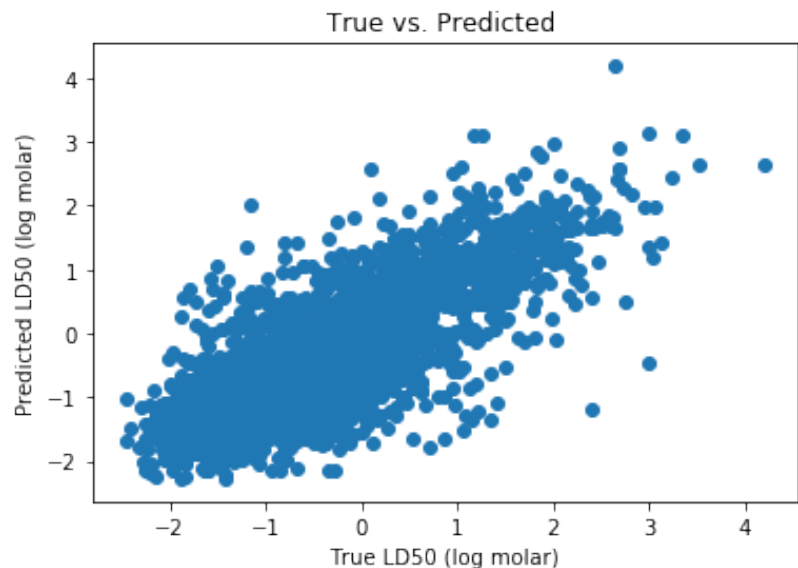
Contact: shah.imran@epa.gov

GenRA - Current research

- Consideration of other information to define and refine the analogue selection & evaluation
 - physicochemical similarity (Helman et al 2018)
 - metabolic similarity (Patlewicz *in prep*),
 - reactivity similarity (Nelms et al 2018)
 - transcriptomics similarity (Tate et al, *under review*)*
- Transitioning to quantitative predictions of toxicity
 - Using GenRA to predict LOAEL, acute oral LD50 (Helman et al 2019a,b)
- Developing a compendium of expert driven read-across examples to investigate how data driven read-across with NAM data can mirror expert assessments (Jenkins et al *in prep*)*

Acute oral toxicity : 'Global' performance

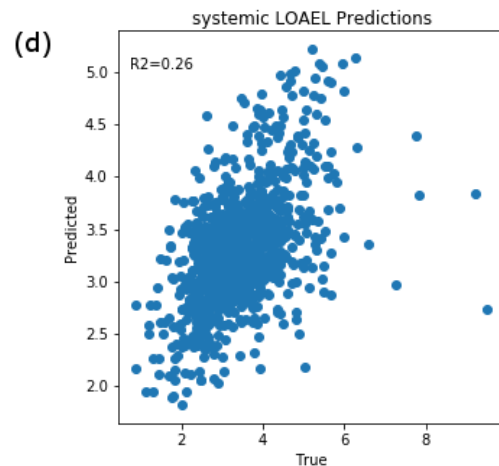
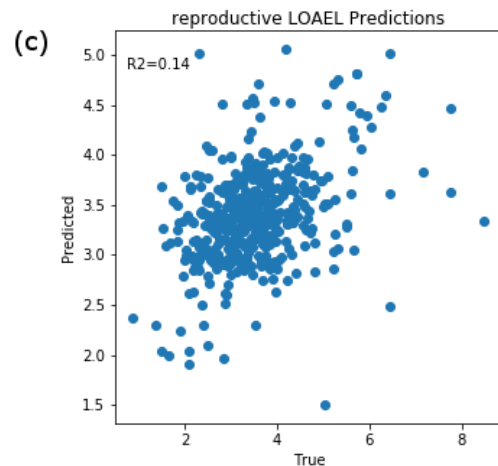
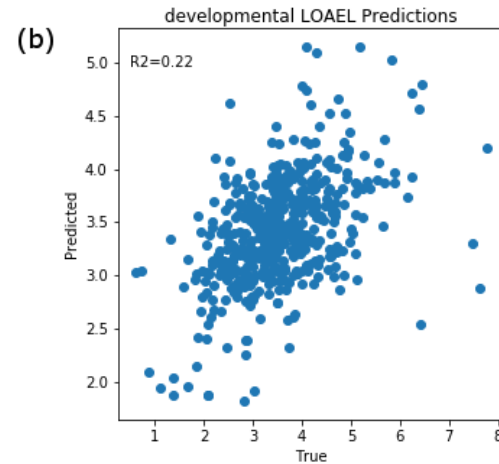
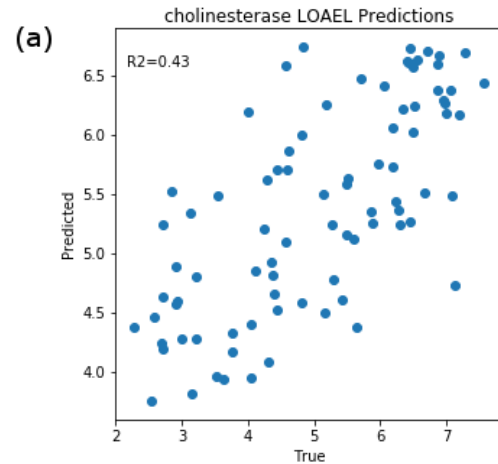
- Search for a maximum of 10 nearest neighbours on entire dataset on the basis of Morgan chemical fingerprints
- Use a min similarity threshold of 0.5



- Linear regression used to fit predicted and observed LD50 values
- $R^2 = 0.61$
- $RMSE = 0.58$

- Monte Carlo CV
- Estimate confidence in R^2
- 75-25 train-test splits
- R^2 values range from 0.46 to 0.62

LOAEL prediction : 'Global' performance



GenRA Predictions using Morgan fingerprints with $k=10$ and $s=0.05$ (mean aggregated LOAELs)
Linear regression used to fit predicted and observed LOAEL values

Endpoint Category	R2
Cholinesterase	0.43
Developmental	0.22
Reproductive	0.14
Systemic	0.26

Characterising metabolic similarity

	Parent_DTXSID	Frag	Parent_smiles	Metabolite_smiles						
0	DTXSID20375106	[#6](=[#8])(-[#8])-[#6]>>[#6]	O=C(O)C(F)(F)OC(F)(F)C(F)(F)OC(F)(F)C(=O)O	O=C(O)C(F)(F)OC(F)(F)C(F)(F)OC(F)F						
1	DTXSID7027831	[#6]-[#7]>>[#7]	CN(CCO)S(=O)(=O)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(...	O=S(=O)(NCCO)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(...						
2	DTXSID7027831	[#6]>>[#8]=[#6]	CN(CCO)S(=O)(=O)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(...	CN(CC(=O)O)S(=O)(=O)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(...						
3	DTXSID7027831	[#6](-[#6].	metab_fp_0	metab_fp_1	metab_fp_2	metab_fp_3	metab_fp_4	metab_fp_5	metab_fp_6	F)C(F)(F)C(F)(F)C(F)(F)...
		DTXSID00190950	0	0	0	0	0	0	0	
4	DTXSID8051419	[#6]-[#7+]	DTXSID00192353	0	0	0	0	0	0	C(F)(F)C(F)(F)C(F)(F)
		[#8]-[#6]	DTXSID00194615	0	0	0	0	0	1	0
			DTXSID00379268	0						
			DTXSID00379884	0						
					DTXSID00190950	DTXSID00192353	DTXSID00194615	DTXSID00379268	DTXSID00379884	DTXSID00379884
				DTXSID00190950	1.0	0.0	0.0	0.0	0.0	0.0
				DTXSID00192353	0.0	1.0	0.0	0.5	0.0	0.0
							1.0	0.0	0.0	0.0
							0.0	1.0	0.0	0.0
				DTXSID00379884	0.0	0.0	0.0	0.0	1.0	1.0

Creating custom fingerprints to characterise metabolic transformations

Creating custom fingerprints to
characterise metabolic transformations

GenRA - Overall goal

- Quantify the contribution that different similarity contexts play in toxicity prediction and how that differs depending on the toxicity endpoint of interest, the chemical of interest and whether it mirrors expert driven read-across
 - Quantify level of confidence for prediction made
- => objective, reproducible read-across assessments

GenRA Summary

- GenRA is an attempt to move towards an objective read-across approach where uncertainties and performance can be quantified. Provides opportunities for NAM data to be incorporated.
- GenRA v1.0 establishes a baseline in performance. The approach relies on chemical descriptors to predict binary toxicity values but work continues to characterise other contexts of similarity (e.g. mechanistic, reactivity, metabolism) and quantify their contribution in predicting *in vivo* toxicity outcomes.
- GenRA v1.0 exists as an app within the Dashboard to facilitate a workflow approach to make read-across predictions. An updated version is anticipated this summer. A python package (genra-py) has been released (March 2021) to facilitate batch processing using user specific datasets.
- Items* will be presented at QSAR2021 - see qsar2021.org

ICCVAM Read-Across Workgroup

- In 2018, US Agencies established a read-across workgroup (RAWG) under ICCVAM to develop and implement a plan to build capacity in the development and application of read-across approaches and to harmonise them.
- Initially, the RAWG summarised current experiences and needs, and catalogued the different tools applied (Patlewicz et al (2019))
- More recent RAWG efforts have been focused on developing a compendium of member agency read-across case studies to inform guiding principles for different read-across decision contexts.
- Several case studies were discussed ranging from the utility of metabolic data in categories, the FDA Extended Decision Tree for TTC to the qualitative use of ToxCast data to characterise bioactivity similarity of a target and candidate analogues.
- A short manuscript is in preparation to summarise the case studies and extract any general guiding principles that will be informative as part of ongoing efforts to refine existing guidance e.g. OECD grouping guidance.