# Dataset Curation, Assessment of their Quality, and Prediction Model Developments for Safe and Sustainable Nanotechnology (S2NANO)

Prof. Tae Hyun Yoon, CEO/Ph.D.

Hanyang University & Yoon Idea Lab. Co. Ltd.



### References - S<sup>2</sup>NANO:PredictNano

J.S. Choi, T.X. Trinh, <u>T. H. Yoon</u>, J. Kim<sup>\*</sup>, H. G. Byun<sup>\*</sup> (2019) <sup>"</sup>Quasi-QSAR for predicting the cell viability of human lung and skin cells exposed to different metal oxide nanomaterials "*Chemosphere*, 217, 243

J.S. Choi, M.K. Ha, T.X. Trinh, <u>T.H.Yoon</u>, H.G.Byun\* (2018) **"Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources" Scientific Reports**, 8, 6110

M.K. Ha, T.X. Trinh, J.S. Choi, D. Maulina, H.G. Byun<u>, T.H.Yoon</u>\*, (2018) **"Toxicity Classification of Oxide Nanomaterials:** *Effects of Data Gap Filling and PChem Score-based Screening Approaches.* "*Scientific Reports*, **8**, 3141

T.X. Trinh, M.K. Ha, J.S. Choi, H.G. Byun, <u>T.H. Yoon</u>\* (2018) **"Dataset Curation, Assessment of their Quality and Completeness, and nanoSAR Classification Model Development for Metallic Nanoparticles**" **Environmental Science: Nano** 5, 1902-1910

*T.X. Trinh, J.S. Choi, H. Jeon, H.G. Byun, <u>T.H.Yoon</u>\*, J.Kim\*, (2018) "Quasi-SMILES based Nano-QSAR model to predict the cytotoxicity of multi-walled carbon nanotubes to human lung cells." <i>Chemical Research in Toxicolology*, 31(3), 183

Sunil Kr Jha\*, <u>TH Yoon</u>, Zhaoqing Pan (2018) "Multivariate statistical analysis for selecting optimal descriptors in the toxicity modeling of nanomaterials", **Computers in Biology and Medicine**, 99, 161

*D.W. Boukhvalov\*, <u>T.H.Yoon</u>, (2017) "Development of Theoretical Descriptors for Cytotoxicity Evaluation of Metallic Nanoparticles." <i>Chemical Research in Toxicolology*, 30 (8), 1549

# How can we bridge the gap between Excellence in nano R&D and Profitable nano Industry ?





<u>Safe & Sustainable Nano</u>technology Portal : A Toolbox for Nanomaterial Characterization & Safety Assessments www.s2nano.org



# Measurements to Models



### **User Friendly Nanosafety Prediction System**



Characteristics of Nanosafety Data ?



Dataset curation based on the assessment of data quality / completeness

Development of generalized prediction models with wider applicability domains

## Model Development Workflow in S<sup>2</sup>NANO

Core Dataset from our Own Experiments Extended Dataset from Literature Mining

- Info.DB, Mat.DB, QM DB,
- PChem. DB, Tox DB (in vitro, in vivo, eco).,





### **Data Collection - Experimental**



Core Dataset

from our Own Experiments

- Core : Tox(in vitro) DB



## **Data Collection - Literature Mining**



### Extended dataset

- from Literature Mining
- Info.DB, Mat.DB, QM DB,
- PChem DB, Tox DB (in vitro, in vivo, eco).,



Core & Extended dataset

from Experiment & Literature Mining

- Info.DB, Mat.DB, QM DB,
- PChem DB, Tox DB (in vitro, in vivo, eco).,



# **Collected Database**



Info ID     Authors     Manufacturer     AMA175															- 24	
- Authors - Manufacturer					Title											소기와
Manufacturer				•	Keyword											
, AHAI TL				· ·	Product #											
-00/1					Journal											
- score	kkkkk.	*****														
																2078
Validation	c total score 👻	score	T Score	생성일	생성자	수정일	수정자	Info #	Source Type	Titles		DOI number	Authors	Affiliation	n Fun	nding Ag
1 미겸증	*****	****	****	2015-02	T021	2015-05		23075393	article	Determining the	pharm acokine	··· 10.1021/nl3···	Malfat···	Bioscienc	ces ···	
2 검증 완료	****	*****	****	2013-08	박세묘	2015-05…	윤태현	22112499	article	Zinc oxide nano	particles interfe	10.1093/toxs	· Kao Y…	Departme	ent ···	
3 미겸증	****	*****	****	2016-05	- 김수진	2016-08…	김수진	24894644	article	The performanc	e of gradient al	···· 10.1016/j.bi··	Soen…	MoSAIC/8	'Bio ···	
4 미검증	*****	*****	+++++	2015-02-	70.04	DOLE OF		00000750	a shind a	And an and from the form		10 1001 1-0				
Material DI     Validation che     Info ID     Material group     Mauntacturer	B cks ⊚ 검원	) 완료 🔘 미검종 (	) 수정묘청		Source Type Mat # Material type	선택 .	¥	20020133	annure	Interaction of go	old nanoparticle	10.1021/ma-	Lacer	Center fo	or Bi···· 김백	초기화
• Material DI • Validation cher • Info ID • Material group • Manufacturer • 생성자	B cks @ 검 문	ੇ 완료 ○ 미검종 (	) 수정요청		Source Type Mat # Material type Product #	전택	▼ 	20020193	ann	interaction of go	old nanoparticle	10.1021/ma-	Lacer	Center fo	ar Bi… 김백	초기화
• Validation che - Info ID - Material group - Mand'acturer - 생성자 - score	B ricks 이 검절	· ***	<ul> <li>● 수정요청</li> </ul>		Source Type Mat # Material type Product #	2015-05····	▼ 	20020193	ann	interaction of go	old nanoparticle	UL 1021/mm	Lacer	Center fo	or Bi···· 김백	초기화
• Material DI • Validation che • Info ID • Material group • Manufacturer • 생성자 • score	B rcks @ 검절	: 환료 ) 미검종 (	) 수정요청		Source Type Mat # Material type Product #	전택 ·	¥.	20020193		interaction of go	old nanoparticle		Lacer	Center fo	24백	초기화
• Material DI • Validation che • Info ID • Material group • Manufactuer • 상업자 • score	B cks 22 cks 22 cks cks cks cks cks cks cks cks cks	: 환료 이 미검종 ( score	· 수정요청	·····································	TU21 Source Type Mat # Material type Product #	2015-05-07	v Info #	20020193	Source Type	e Material grou	Material type	Product # Lot #	Mar	Center fo	28백 전쟁	초기화 20개 Tink ta
<ul> <li>Material DI</li> <li>Validation che</li> <li>Info ID</li> <li>Material group</li> <li>Mandactuar</li> <li>생성자</li> <li>Score</li> <li>Validation</li> <li>검증 완료</li> </ul>	B B C T Score → ★★★★★	: 환료 이 미검종 ( - *******	· 수정요청 · 수정요청 · · · · · · · · · · · · · · · · · · ·	·····································	Source Type Mat # Material type Product # 2015-10	2015-05 ····	v Info # 22303956	Mat# Mat1	Source Type article	<ul> <li>Material grou</li> <li>Oxdes</li> </ul>	Material type 1 Ti02	Product # Lot #	Mar	Center fo	2백 전백	초기화 20개
· Material Di           · Validation che           · Info ID           · Naterial group           · Mand actuer           · edata           · eccre           Validation           1 검증 완료           2 검증 완료	B icks		· 수정요청 생정일 2014-03 2014-03		Source Type           Mat #           Material type           Product #           2015-10···           2015-10···	2015-05····	▼ Info # 22303956 22303956	Mat # Mat 4	Source Type article article	Material grou Oxides Oxides	Material type TiO2 SiO2	Product # Lot #	Mai	Center fo	2백 전백	초기화 20개 r link to
· Validation che · Validation che · Info ID · Mand'acturer · 생성자 · Score · Validation · Mand'acturer · 생성자 · Score · Validation · Mand'acturer · 생성자 · Score · Validation che · National group · Mand'acturer · Walidation · Mand'acturer · Walidation che · Mand'acturer · Walidation · Mand'acturer · Walidation	B cks ○ 22 ☆☆☆☆☆ ~ ( T Score ~ ☆☆☆☆☆☆	<ul> <li>★★★★★</li> <li>★★★★★</li> <li>★★★★★</li> <li>★★★★★</li> <li>★★★★★</li> </ul>	· 수정도성 생상일 2014-03… 2014-03… 2014-03…	생성자 박종훈 박종훈	Source Type           Mat #           Material type           Product #           2015-10···           2015-10···           2015-10···	2015-05-55	▼ Info # 22303956 22303956	Mat # Mat I Mat 1 Mat 3	Source Typp article article article	Material grou Oxdes Oxdes	Material type 1 Tro2 SiO2 ZnO SiO2 ZnO SiO2 SiO2 SiO2 SiO2 SiO2 SiO2 SiO2 SiO	Product # Lot #	Mar	Center fo	<mark>감박</mark> Manufacturin	초기화 20개

																2074 •
Validation ch	T Score	score	생성일	생성자	수정일	수정자	Info #	Mat#	Manufacture	Product #	PC #	Source Type	pchem type	Core (or Grai	Internal Dian	External [
미검증	*****	*****	2016-08…	김수진	2016-10	김수진	19041333	Mat 1	Carbon N…		PC 1	article	pchem	1 ± 0.2 nm		
2 미겸증	****	****	2016-08	김수진	2016-09	김수진	20016928	Mat 1	Nano Lab …		PC 1	article	pchem	$22.5~\pm~7.\cdots$		
미검증	*****	****	2016-08	김수진	2016-08	김수진	21651974	Mat 1	Sigm a Al…	704113	PC 1	article	pchem	1 ± 0.3 nm		
1 미검증	*****	*****	2016-08	김수진	2016-10	김수진	21651974	Mat 2	Sigma Al…	652490	PC 1	article	pchem	$1.4\pm0.1\cdots$		

TOXICITY DB	- In vitro																	
- Validation chec - Info # - Material group - Manufacturer	ks 0 겸종	8 완료 💿 미검종 🔘	수정묘형		<ul> <li>Source Type</li> <li>Mat #</li> <li>Material type</li> <li>Product #</li> </ul>	선택	T									겁	백 초기	화
생성자 score	hakakak-	*****																
																	2	2078 1
Validation	(T Score +	score	생성일	생성자	수정일	수정자	info #	Mat #	Manufacture	Product #	PC #	Tox #	Source Typ	Туре	Value	Assay meti	2 Cell line Na	07H Spec
Validation 1 검증 완료	<pre>T Score → ★★★★★★</pre>	score ★★★★★	생성일 2014-0…	<b>생성자</b> 김수진	수정일 2016-1…	<b>수정자</b> 김진배	Info # 22303956	Mat# Mat 4	Manufacture	Product #	PC#	Tox #	Source Typ article	Туре	Value	Assay meti	2 Cell line Na MCF7	07H Spec Huma
Validation           1         검증 완료           2         검증 완료	T Score ~       *****	score	생성일 2014-0… 2014-0…	<b>생성자</b> 김수진 김수진	수정일 2016-1… 2016-1…	<b>수정자</b> 김진배 김진배	Info# 22303956 22303956	Mat# Mat 4 Mat 3	Manufacture	Product #	PC # PC 2 PC 1	Tox # Tox 1 Tox 1	Source Type article article	Туре	Value	Assay metil MTT MTT	2 Cell line Na MCF7 MCF7	Spec Hum Hum
Validation           1         검증 완료           2         검증 완료           3         검증 완료	( T Score → ★★★★★ ★★★★★ ★★★★★★	score ***** *****	생성일 2014-0… 2014-0…	생성자 김수진 김수진 김수진	수정일 2016-1… 2016-1… 2016-1…	<b>수정자</b> 김진배 김진배 김진배	Info # 22303956 22303956 22303956	Mat# Mat 4 Mat 3 Mat 3	Manufacture	Product #	PC# PC2 PC1 PC2	Tox # Tox 1 Tox 1 Tox 1 Tox 1	Source Typ article article article	Туре	Value	Assay metil MTT MTT MTT	2 Cell line Na MCF7 MCF7 MCF7	Spec Huma Huma Huma

### Improve Data Quality via Scoring Methods

### Model Development Workflow in S<sup>2</sup>NANO



Data Completeness (Missing Data Problem) Data Quality (Heterogeneous Source of Data) Data Imbalance (NonToxic >> Toxic )



## Data Preprocessing - Original Dataset (Dataset I)

- Total 20 attributes, but only 14 attributes were used as Descriptors
  - Dose(1) / Pchem(8) / Tox(7) / QM(4) attributes
  - Measurement Methods attributes were not used for Model development (-4)
  - Cell name attribute was not used for Model development (-1)
  - Cell Viability was used for Toxic/Non-Toxic Endpoint (-1)

- Toxic/Non-Toxic as Endoint
  - Toxic when Cell Viability < 50 %
  - NonToxic when Cell Viability  $\geq$  50 %

	А	В	F	G	н	1		J	К	L	М	Ν	0	Р	Q
1	Pubmed ID N	Material type	Method core size	Hydro size (nm)M	lethod hydro size	Surface charge (i	mV) Method s	urface charge	Surface area (m2/g)	Method surface area	∆Hsf (eV)	Ec (eV)	Ev (eV)	χMeO (eV)	Assay
2	22502734 A	AI2O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
3	22502734 A	AI2O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	- <mark>9.81</mark>	5.67 MTS	
4	22502734 A	412O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
5	22502734 A	AI2O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
6	22502734 A	4203	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
7	22502734 A	AI2O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
8	22502734 A	412O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
9	22502734 A	412O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
10	22502734 A	412O3	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
11	22502734 A	4203	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 MTS	
12	22502734 A	4203	TEM	260.4 D	LS		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 ATP	
13	22502734 A	412O3	IEM	260.4[D	LS .		0 Zeta-poten	tial			-17.345	-1.51	-9.81	5.67 ATP	
813	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
814	22502734 Z	ZrO2	TEM	312.3 D	LS		12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
815	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
816	22502734 Z	ZrO2	TEM	312.3 D	LS		12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
817	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
818	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
6819	22502734 Z	ZrO2	TEM	312.3 D	LS		12.8 Zeta-poten	tial	i 		-11.252	-3.19	-8.23	5.62 ATP	
820	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 ATP	
821	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
822	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
823	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
824	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
825	22502734 Z	2rO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
826	22502734 Z	ZrO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
827	22502734 Z	2rO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
828	22502734 2	2rO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
829	22502734 Z	2rO2	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	
830	22502734 2	202	TEM	312.3 D	LS	-	12.8 Zeta-poten	tial			-11.252	-3.19	-8.23	5.62 LDH	01
831	24983896 Z	2r02	TEM	661.4 D	LS		-8.5 Zeta-poten	tial			-27.67	-1.03	-10.92	4.95 Cyto I o	X-GI0
832	24983896 Z	2rO2	TEM	661.4 D	LS		-8.5 Zeta-poten	tial			-27.07	-1.03	-10.92	4.95 Cyto I o	X-GIO
0033	24983890 Z	7-02		001.4 D	LS		-8.5 Zeta-poten	tial			-21.01	-1.03	-10.92	4.95 Cylo10	x-Glo
0025	24983896 2	7:02		001.4 D			-8.5 Zeta-poten	tiol			-27.07	-1.03	-10.92	4.95 CyloTo	x-Glo
000	24903090 2	702		001.4 D	10		-0.5 Zeta-poten	tiol			-27.07	-1.03	-10.92	4.95 CytoTo	x-Glo
000	24903090 2	7:02		001.4 D	10		-0.5 Zeta-poten	tiol			-27.07	-1.03	-10.92	4.95 Cyl010	x-Glo
0001	24903090 2	702		001.4 D	10		-0.5 Zeta-poten	tiol			-27.07	-1.03	-10.92	4.95 CytoTo	x-Glo
020	24903090 2	7:02		505 D	19		22 Zota poten	tial			-27.07	-1.03	-10.92	4.95 Cylo 10	x-0i0
940	21000333 Z	7rO2		505 D	19		23 Zeta poten	tial			-11.202	-3.19	-0.23	5.62 MTT	
9/1	21800953 2	7rO2		505 D	15		23 Zeta poten	tial			-11.252	-3.19	-0.23	5.62 MTT	
842	21800953 2	7rO2		505 D	15		-23 Zeta-poten	tial			-11.202	-3.19	-0.23	5.02 WIT	
9/2	21800953 2	7rO2		505 D	19		23 Zeta poten	tial			-11.252	-3.19	-0.23	5.02 MTT	
844	210000002			505 D	20		20 Zeta-poten	iliui			-11.2.02	-0.10	-0.20	0.02 mil 1	
														i	

- 216 articles selected from ~600 pdf files
- 26 oxide NPs
- 6,842 data rows

- Missing Data in Oxide NPs' Original Dataset (Dataset I)
  - 18 % of Core Size Data
  - 39 % of Hydrodynamic Size Data
  - 41 % of Surface Charge Data
  - 74 % of Specific Surface Area Data



Quality & Completeness Assessment, Data Gap Filling and PChem score based Screening

# Data Gap Filling Method – Nano Read-Across

# Missing data replacement

## **Conventional Approach**

- substitute missing values with mean values of non-missing values

## **Nano Read Across**

- estimation from other properties of the same nanomaterials (e.g., estimating specific surface area from core size).

- information form manufacturer's specification sheet or other references using the same nanomaterials





Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources

Jang-Sik Choi<sup>1</sup>, My Kieu Ha², Tung Xuan Trinh², Tae Hyun Yoon@² & Hyung-Gi Byun@<sup>1</sup>

Choi et al (2018)

# Data Imbalance Issue

□ SMOTE (Synthetic minority over-sampling technique)



Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, *16*, 321-357.

# Toxic 16% : Nontoxic 84 %

SMOTE (Synthetic Minority Over-sampling TEchnique ) ID (Imbalanced Data) vs. BD (Balanced Data)

# **Curated Datasets**

Safe NANO धर्षकेय प्रकृगहरू ustainable पश्च Prec	dictNA	NC	) > [	Datasets	ogout						
About S <sup>2</sup> NANO FindNANO NANOinfo NANOanalysis PredictN		ʻiki	NANOban	k News/Update							
Chemical Research in Toxic Models Quasi-SMILES-Based Nano-Quantitative Structu	M IANO" by Ire-Act	Dat	aset	:S			HOME	> <u>PredictNANO</u> > Datasets			
Relationship Model to Predict the Cytotoxicity o	of _										
Multiwalled Carbon Nanotubes to Human Lung	Cells	· Mate	rial Grou	р		Material Typ	e (ex: ag, au, s	io2)			
		Sele	ect		•						
	<b>)</b>	·Data	set Name	9		Properties					
	.0					None Material type Core size	size		🛿 Dataset Info ক্রেই		
					ŀ	Hydrodynamic	size	•	Dataset ID	D \$0005	
					Search	Reset			Dataset Name	2016_MeOx_KNU	
					Search	Reset			Material Group	Metal Oxide	
		0	ownload	Dataset Name	Material Group	Property	Cou Pchem Sco	re Data Gap Filling Method	Material	Al203, Cu0, Fe203, Fe304, Si02, Ti02, Zn0, Ce0	2, Co3O4, CoO, Cr2O3, Gd2O3, HfO2, In2O3, La2O3, Mn2O3, Ni2O3, NiO, Sb2O3
		1	Ŧ	2017_MeOx_II_KNU	Metal & Metal Oxid	e 8			Number Of Rows	1738 Motorial tuno	
		2	Ŧ	2017_MeOx_I_KNU	Metal Oxide	18	4.7 ± 0.3	Manufacturer's specification Estimation PChem score >=4		Viaterial type Core size Hydrodynamic size Surface charge Surface area	
		з	Ŧ	2017_MWCNT_HYU	Carbon	8				dHsf Ec	
		4	Ŧ	2017_Metal_HYU	Metal	14	4.3 ± 0.3	Manufacturer's specification	Property (18)	EV xMeO Mass do se	
		5	Ŧ	2017_Metal_validation	Metal	13				Exposure tim e Assay m ethod	Details on Dataset
		6	Ŧ	2017_MeOx_validation	Metal Oxide	13				Cell line Cell species	
		7	¥	2017_MeOx_IIIB_HYU	Metal Oxide	21	4.9 ± 0.2	Manufacturer's specification Estimation 20% data with top PChem		Cell type Cell viability Toxicity	
		8	Ŧ	2017_MeOx_IIIA_HYU	Metal Oxide	21	4.8 ± 0.1	Manufacturer's specification Estimation 50% data with top PChem	Data Gap Filling Method	Manufacturer's specifications Estimation 50% data with top PChem score	
Download Dataset		-		2017_MeOx_II_HYU	Metal Ovida	22	47 + 0.2	Manufacturaria specification	Fonenii Soone	4.0 ± 0.1	
		10	Ŧ	2017_MeOx_I_HYU	Metal Oxide	21	2.8 ± 1.3	Estimation Mean substitution			
		14 4	Page	1 of 2 🕨 🕅 🍣	1			Displaying 1 - 10 of 16			

### Model Development Workflow in S<sup>2</sup>NANO

Logistic Regression Algorithm Random Forest Algorithm Support Vector Machine Algorithm Backpropagation Algorithm



# Model Development – Algorithm Selection



# Model Development – Validation (Internal & External)

### Model validation



Jang-Sik Choi<sup>1</sup>, My Kieu Ha<sup>2</sup>, Tung Xuan Trinh<sup>2</sup>, Tae Hyun Yoon<sup>®2</sup> & Hyung-Gi Byun<sup>®1</sup>

# SCIENTIFIC REPORTS

Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources

Jang-Sik Choi<sup>1</sup>, My Kieu Ha<sup>2</sup>, Tung Xuan Trinh<sup>2</sup>, Tae Hyun Yoon<sup>1</sup> & Hyung-Gi Byun<sup>1</sup>

Choi et al (2018)

# **Normalization Method**

		True	False	False	True			Balanced	St	andard deviation	
Algorithm	Normalization method	positive	positive	negative	negative	Sensitivity	Specificity	accuracy	Sensitivity	Specificity	Balanced accuracy
	min-max	39	12	16	278	70.91%	95.86%	83.39%	4.84%	1.18%	2.45%
	z-score	39	12	16	278	70.91%	95.86%	83.39%	4.84%	1.58%	2.47%
LK	log	46	11	9	279	83.64%	96.21%	89.92%	4.45%	1.01%	2.33%
	combination	45	8	10	282	81.82%	97.24%	89.53%	3.99%	1.15%	3.99%
	min-max	28	6	27	284	50.91%	97.93%	74.42%	10.07%	0.72%	5.04%
C)////	z-score	29	7	26	283	52.73%	97.59%	75.16%	9.73%	0.35%	4.87%
5111	log	40	5	15	285	72.73%	98.28%	85.50%	5.48%	0.65%	2.65%
	combination	41	5	14	285	74.55%	98.28%	86.41%	5.48%	0.67%	2.73%
	min-max	45	5	10	285	81.82%	98.28%	90.05%	5.71%	0.64%	2.92%
	z-score	44	5	11	285	80.00%	98.28%	89.14%	4.32%	0.67%	2.25%
KF	log	45	5	10	285	81.82%	98.28%	90.05%	5.02%	0.61%	2.49%
	combination	45	5	10	285	81.82%	98.28%	90.05%	4.93%	0.61%	2.44%
	min-max	38	15	17	275	69.09%	94.83%	81.96%	15.31%	1.07%	7.79%
	z-score	40	6	15	284	72.73%	97.93%	85.33%	5.89%	0.60%	2.98%
ANN	log	43	8	12	282	78.18%	97.24%	87.71%	6.00%	0.64%	3.01%
	combination	48	8	7	282	<mark>87.27%</mark>	97.24%	<mark>92.26%</mark>	4.05%	0.66%	2.11%

# SCIENTIFIC REPORTS

Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources

Jang-Sik Choi<sup>1</sup>, My Kieu Ha², Tung Xuan Trinh², Tae Hyun Yoon@² & Hyung-Gi Byun@<sup>1</sup>

Choi et al (2018)

Data Imbalance Issue

2.33%

7.25%

2.73%

0.56%

2.44%

0.63%

2.11%

0.75%

Normalization	A I	Data	True	False	False	True	O a se a iti a ita a		Balanced	St	andard deviation	
method	Algorithm	Data	positive	positive	negative	negative	Sensitivity	Specificity	accuracy	Sensitivity	Specificity	Balanced accuracy
Log	I D	ID	46	11	9	279	83.64%	96.21%	89.92%	4.45%	1.01%	2.
LUG	LK	BD	243	27	18	234	93.10%	89.66%	91.38%	5.84%	9.31%	7.
Combination	S/M	ID	41	5	14	285	74.55%	98.28%	86.41%	5.48%	0.67%	2.
Compination	3 1 1 1	BD	255	7	6	254	97.70%	97.32%	97.51%	1.15%	0.98%	0.
Combination	DE	ID	45	5	10	285	81.82%	98.28%	90.05%	4.93%	0.61%	2.
Compination	ΓΓ	BD	253	4	8	257	96.93%	98.47%	97.70%	0.88%	0.72%	0.
Combination	ΔΝΙΝΙ	ID	48	8	7	282	87.27%	97.24%	92.26%	4.05%	0.66%	2.
Combination		BD	258	5	3	256	<mark>98.85%</mark>	98.08%	<mark>98.47%</mark>	0.79%	0.85%	0.

# Toxic 16% : Nontoxic 84 %

SMOTE (Synthetic Minority Over-sampling TEchnique ) ID (Imbalanced Data) vs. BD (Balanced Data)

### Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources

Jang-Sik Choi<sup>1</sup>, My Kieu Ha<sup>2</sup>, Tung Xuan Trinh<sup>2</sup>, Tae Hyun Yoon<sup>1</sup> & Hyung-Gi Byun<sup>1</sup>

SCIENTIFIC REPORTS

# **Internal vs. External Validations**

Choi et al (2018)

	Normalization		_	True	False	False	True	<b>_</b>		Balanced	St	andard deviation	
	method	Algorithm	Data	positive	positive	negative	negative	Sensitivity	Specificity	accuracy	Sensitivity	Specificity	Balanced accuracy
	Log	IR	ID	46	11	9	279	83.64%	96.21%	89.92%	4.45%	1.01%	2.33%
Inte	rnal valida	ntion	BD	243	27	18	234	93.10%	89.66%	91.38%	5.84%	9.31%	7.25%
iiiic		SVM	ID	41	5	14	285	74.55%	98.28%	86.41%	5.48%	0.67%	2.73%
	Combination	3 1 10	BD	255	7	6	254	97.70%	97.32%	97.51%	1.15%	0.98%	0.56%
	Combination	DE	ID	45	5	10	285	81.82%	98.28%	90.05%	4.93%	0.61%	2.44%
	Combination		BD	253	4	8	257	96.93%	98.47%	97.70%	0.88%	0.72%	0.63%
	Combination	ΔΝΙΝΙ	ID	48	8	7	282	87.27%	97.24%	92.26%	4.05%	0.66%	2.11%
	Combination		BD	258	5	3	256	<mark>98.85%</mark>	98.08%	<mark>98.47%</mark>	0.79%	0.85%	0.75%
	Normalization			Truo	Falso	Falso	True			Balanced	St	andard deviation	

	Normalization			True	False	False	True	<b>•</b> • • • •	<b>a</b> 10 1	Balanced	3		
	method	Algorithm	Data	positive	positive	negative	negative	Sensitivity	Specificity	accuracy	Sensitivity	Specificity	Balanced accuracy
	Log	ID	ID	25	6	4	194	86.21%	97.00%	91.60%	7.95%	0.70%	3.86%
Evto	rnalvalida		BD	26	21	3	179	89.66%	89.50%	89.58%	9.10%	7.49%	5.46%
LYIG	Combination		ID	22	5	7	195	75.86%	97.50%	86.68%	11.64%	1.06%	5.71%
	Combination	2010	BD	25	10	4	190	86.21%	95.00%	90.60%	8.85%	1.85%	3.79%
	Combination	DE	ID	24	3	5	197	82.76%	98.50%	90.63%	12.75%	1.12%	6.23%
	Combination	KF	BD	25	9	4	191	86.21%	95.50%	90.85%	11.24%	1.84%	5.84%
	Combination		ID	23	4	6	196	79.31%	98.00%	88.66%	8.74%	1.42%	4.42%
	Combination	AININ	BD	27	13	2	187	<mark>93.10%</mark>	93.50%	<mark>93.30%</mark>	6.41%	0.83%	3.08%



# Performance Comparisons of nanoSAR classification Models

Mardal Nama	Dataset	A las a state as	Internal validation	<b>External validation</b>	Dublications
		Algorithm	Accuracy or R <sup>2</sup>	Accuracy or R <sup>2</sup>	Publications
PM-100	2017_Metal_HYU	Random forest	86%	82%	
PM-101	2016_Metal_KNU	Backpropagation	95%	67%	
PM-102	2015_Metal_HYU	Support vector machine	90%	75%	
PM-103	2015_Metal_MeOx_KNU	Backpropagation	91%	86%	
PM-104	2015_Metal_HYU	Support vector machine	90%	82%	
PM-106	2017_MeOx_I_KNU	Resilient backpropagation	93%	86%	
PM-107	2017_MeOx_II_HYU	Random forest	94%	54%	Ha et al. (2018) Scientific Reports
PM-108	2016_MeOx_KNU	Backpropagation	95%	85%	Choi et al. (2018) Scientific Reports
PM-114	2017_MeOx_I_KNU	Random forest	91%	87%	
PM-115	2017_MeOx_I_KNU	Support vector machine	93%	89%	
PM-116	2017_MWCNT_HYU	Quasi-QSAR	0.89	N/A	Trinh et al. (2018) Chem. Res. In Toxicology
PM-117	2017_MeOx_II_KNU	Quasi-QSAR	0.79	0.71	Choi et al. (2018) Chemosphere
PM-118	2015_Metal_KNU	Backpropagation	88%	67%	
PM-120	2017_Metal_HYU	Support vector machine	86%	82%	Trinh et al. (2018) Environmental Science : NANO



Towards a generalized toxicity prediction model for oxide nanomaterials using integrated data from different sources

Jang-Sik Choi<sup>1</sup>, My Kieu Ha<sup>2</sup>, Tung Xuan Trinh<sup>2</sup>, Tae Hyun Yoon<sup>2</sup> & Hyung-Gi Byun<sup>1</sup>

Choi et al (2018) (Impact Factor = 4.259)

### **Environmental** Trinh et al (2018) (Impact Factor = 6.047) Science Nano Curation of datasets, assessment of their quality and completeness, and nanoSAR classification model development for metallic nanoparticles<sup>†</sup> Tung X. Trinh, <sup>(i)</sup><sup>a</sup> My Kieu Ha,<sup>a</sup> Jang Sik Choi,<sup>b</sup>

Tung X. Trinh, 💷 a My Kieu Ha, a Jang Sik Cho Hyung Gi Byun 🗊 b and Tae Hyun Yoon 🗊 \*a

# **Prediction Models**



Ma	teria	l Group			· Mat	terial Type	(ex: ag, au	, sio2)	
S	elect	:		Ŧ					
Mo	odel I	Name			· Pro	perty			
					None Mate Core Hydr	e erial type size odynamic	size		•
				Search	h Res	et			
	DW	Model Name		1	Material Group	Accuracy	Sensitivity	Specificity	End Point
11	☑	2011 Small Ron	g Liu	1	Metal Oxide	100	0	0	Cytotoxicity (to
12	Z	. 016 MetalOxic	le BP KNU	1	Metal Oxide	95	86	87	Cytotoxicity (to.
13		2017 MetalOxic	le RF HYU	1	Metal Oxide	94	65	98	Cytotoxicity (to.
14	2	2017 MetalOxid	le BP KNU		Metal Oxide	93	90	94	Cytotoxicity (to.
15	2	Nano-Lazar			Metal	0	0	0	
16		2015 Metal_Me	talOxide SVM HYU		Metal & M	90	54	97	Cytotoxicity (to.
17	2	2015 Metal_Me	talOxide BP KNU		Metal & M	91	48	97	Cytotoxicity (to.
18	2	2015 Metal SVN	и нуџ		letal	89	47	96	Cytotoxicity (to.
19		2016 Metal BP	KNU	1	M. tal	95	79	96	Cytotoxicity (to.
20	ø	2017 Metal RF	HYU	1	Metal	86	26	99	toxic / Nontoxi.
4	<b>∢</b>	Page 2 of	2   🕨 🕅   🧞					Dis	playing 11 - 20 of
	<b>N</b>	1odel Info 조회							
	Mo	del ID	PM-108						
	Mo	del Nam e	2016 MetalOxide BP K	NU					
	Ma	terial Group	Metal Oxide						
	Als	iorithm	Backpropagation						
	Ors	ganization	KNU	Specificity = 00	52% Acouraci -	00 27%			
	Pe	form ance	Accuracy : 95 Sensitivity : 86 Specificity : 97	, opecificity = 98	.55%, ACCURACY =	00.21%			
	En	d Point	Cytotoxicity (toxicity la	bels)					
	Та	rget Cell Line	Diverse	<b>D</b> · · ·				1	
				Detai	ls on	Fach	Pred	dictic	$n M \alpha$

0.5>=

Toxic Range



### Model Development Workflow in S<sup>2</sup>NANO



Relative Importance of Attributes Applicability Domains



Importance of each attribute to toxicity is represented by their weights

Jang-Sik Choi<sup>1</sup>, My Kieu Ha<sup>2</sup>, Tung Xuan Trinh<sup>2</sup>, Tae Hyun Yoon 🎯<sup>2</sup> & Hyung-Gi Byun 🎯



OPEN Toxicity Classification of Oxide Nanomaterials: Effects of Data Gap Filling and PChem Score-based Ha et al.(2018) Screening Approaches (Impact Factor = 4,259)

# **Applicability Domains of Models**

	Ι		II		III-A		III-B	
Attribute	Min.	Max.	Min.	Max.	Min.	Max.	Min.	Max.
Dose (µg/mL)	0	10000	0	167000	0	1500	0	1500
Time (h)	0	360	1	168	2	72	6	72
Core size (nm)	2.7	629	2.7	496	5	496	5.9	193
Hydro. size (nm)	8.6	6181	8.6	2300	12.5	1463	12.5	1457
Surface charge (mV)	-63.3	61.9	-63.3	61.9	-52	61.9	-47.6	42.8
Surface area (m <sup>2</sup> /g)	0.8	1150	5.5	576	5.5	576	6	576
$\Delta H_{sf}(eV)$	-64.7	-1.2	-64.7	-1.2	-26.8	-1.2	-26.8	-1.6
E <sub>c</sub> (eV)	-6.6	-0.1	-6.6	-0.1	-5.2	-0.1	-5.3	-0.3
E <sub>v</sub> (eV)	-11.4	-5.0	-11.3	-5.0	-11.1	-5.0	-11.4	-5.0
χ (eV)	3.2	8.3	3.4	8.3	3.4	6.8	3.8	8.3

 Table 4. Applicability domains regarding the numerical attributes.

Environm	ental Trinh et al (2018)	
Nano	Curation of datasets, assessment of their quality and completeness, and nanoSAR classification model development for metallic nanoparticles†	
	Tung X. Trinh, 💿 <sup>a</sup> My Kieu Ha, <sup>a</sup> Jang Sik Choi, <sup>b</sup> Hyung Gi Byun 🔞 <sup>b</sup> and Tae Hyun Yoon 🔞 * <sup>a</sup>	

Table S6. Applicability domain of nanoSAR models built from datasets A, B and C.

Attribute	Dataset A	Dataset B	Dataset C
NPs type	Ag, Au	Ag, Au	Ag, Au
Shape	Particle, hollow, nanorod	Particle, hollow, nanorod	Particle, hollow, nanorod
Core size (nm)	2 - 120	2 - 120	2.5 - 120
Hydrodynamic size (nm)	7.1 – 300	7.1 – 300	7.1 - 300
Surface charge (mV)	-78.8 - 58.2	-78.8 - 58.2	-78.8 - 58.2
Specific surface area (m <sup>2</sup> /g)	2.3 - 185.7	2.3 - 185.7	2.3
Dose (ppm)	0 - 400	0 - 400	0 - 400
Exposure time (h)	0 – 96	0 – 96	0 – 96



Implementation of Collected Database, Curated Datasets and nanoSAR classification models in S<sup>2</sup>NANO portal.



### SUMMARY

- To overcome current issues in nanosafety data, such as small, unbalanced, & heterogeneous datasets with many missing values, we have collected a comprehensive nanosafety database (S2NANO) from experiments as well as literature mining. (33,393 rows of raw data were collected)
- These data were further processed and 16 quality screened datasets were curated : Data gap-filled with nano read-across methods and assessed their data quality / completeness based on Pchem score. Using these curated datasets, 13 prediction models were developed with different algorithms (LR, SVM, RF, ANN) and validated internally & externally.
- These comprehensive database, curated datasets, and nanosafety prediction models were implemented in S2NANO portal with user-friendly interfaces for future applications in safety by designand regulation compliance.





Measurements & Models for Nanomaterials

### QSAR vs Quasi-QSAR

[Problem] In the case of nanomaterials, the molecular structures of nanomaterials is the same as bulk chemicals

In the case of "classic" QSPR/QSAR analysis the paradigm is the following:

→ Endpoint is a mathematical function of molecular structure



Ref. Toropova, Alla P., and Andrey A. Toropov. "Optimal descriptor as a translator of eclectic information into the prediction of membrane damage by means of various TiO 2 nanoparticles." *Chemosphere* 93.10 (2013): 2650-2655. Ref. Toropov, Andrey A., Robert Rallo, and Alla P. Toropova. "Use of quasi-SMILES and Monte Carlo optimization to develop quantitative feature property/activity relationships (QFPR/QFAR) for nanomaterials." *Current topics in medicinal chemistry* 15.18 (2015): 1837-1844.



# Acknowledgements



# Acknowledgements

# Safe NANO ustaínable http://portal.s2nano.org

Prof. BYUN, Hyung-Gi (KNU) Dr. CHOI, Jang-Sik (KNU/HYU) Mr. CHOI, Woosu (TO21) Mr. KIM, Doyoung (TO21) Dr. BAE, Hee-Kyung (TO21) Dr. KIM, Jongwoon (KIST-Europe/KRICT) Dr. JEON, Hyun Pyo (KIST-Europe) Dr. YOON, Seokju (KIT) Dr. Oh, Jeong-Hwa (KIT) Ms. HA, Kieu My (HYU) Mr. TRINH, Xuan Tung (HYU)

**Collaborators** 



# Thanks for your attention

### Model Interpretation : Attribute Importance



□ Mann-Whitney-Wilcoxon-test (also called the Wilcoxon rank-sum test)

H0 : the distributions are the same H1 : the distributions are not the same

**Galaxie Relative importance** 



Ibrahim, O. M. "A comparison of methods for assessing the relative importance of input variables in artificial neural networks." Journal of Applied Sciences Research 9.11 (2013): 5692-5700.

### Model Interpretation : Applicability Domain

### □ KNN-based applicability domain



The new compound will be predicted by the model, only if:

 $\mathbf{D}_i \leq <\mathbf{D}_k > +\mathbf{Z} \times \mathbf{s}_k$ 

### With Z, an empirical parameter(0.5 by default)

 $< D_k >$ : average Euclidian distance between each compound of the training set and its k nearest neighbors in the descriptors space.  $s_k$ : standard deviation of the distances between each compound of the training set and its k nearest neighbors in the descriptors space.  $D_i$ : the average of the distances between i and its k nearest neighbors in the training set.



Tropsha, Alexander, Paola Gramatica, and Vijay K. Gombar. "The importance of being earnest: validation is the absolute essential for successful application and interpretation of QSPR models." Molecular Informatics 22.1 (2003): 69-77.