

Optimierung und Sicherung der
Datenqualität von STOFF-IDENTV. Gronau¹, Wielenbach/D, M. Luthardt³, Weißenstephan/D, T. Placht³, Weißenstephan/D, A. Gilg³,
Weißenstephan/D, F. Leßke³, Weißenstephan/D, M. Sengl², Augsburg/D, M. Letzel¹, Wielenbach/D¹Bayerisches Landesamt für Umwelt, Demollstraße 31, 82407 Wielenbach²Bayerisches Landesamt für Umwelt, Bürgermeister-Ulrich-Str. 160, 86179 Augsburg³Hochschule Weißenstephan-Triesdorf, Vöttinger Str. 27, 85354 Freising

Bayerisches Landesamt für Umwelt, Demollstraße 31, 82407 Wielenbach, veronika.gronau@lfu.bayern.de

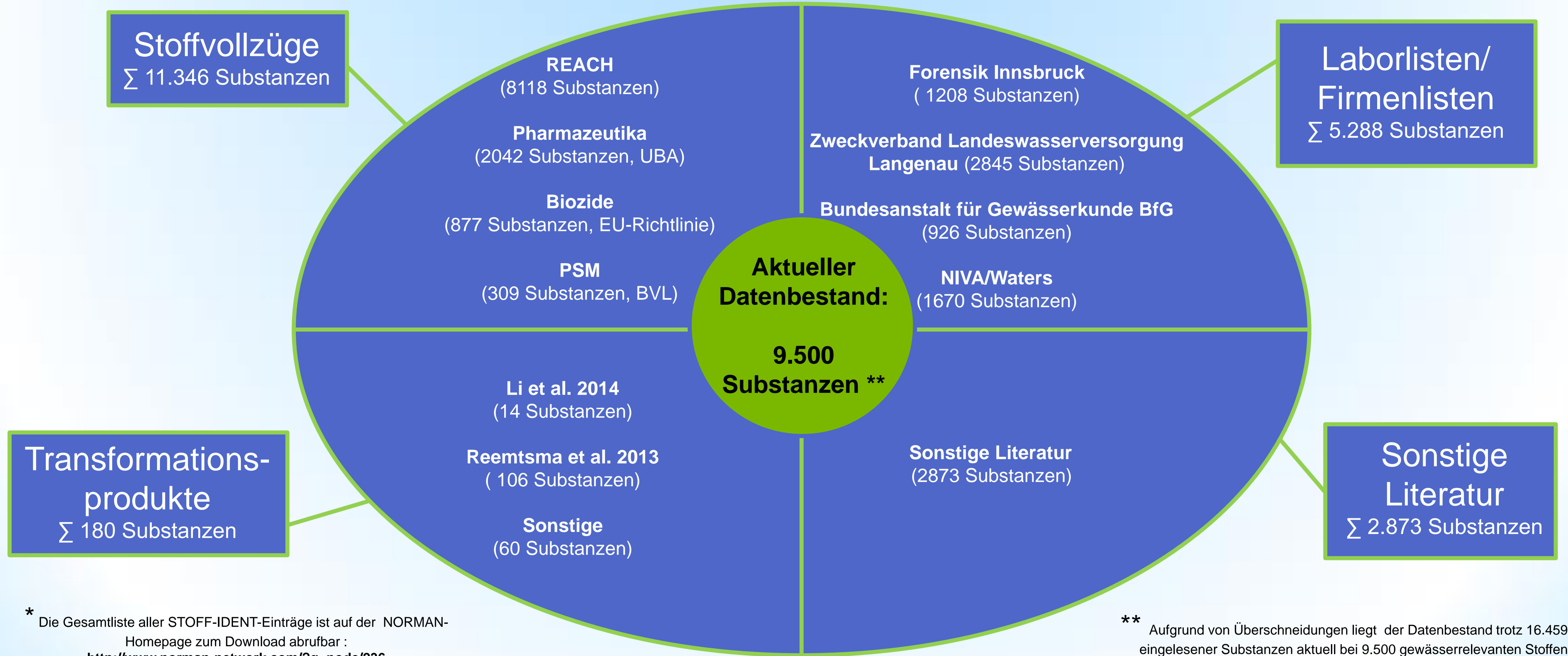
STOFF-IDENT- Datenbank

- Hilfsmittel zur Identifizierung bislang unbekannter gewässerrelevanter Stoffe mithilfe der Non- und Suspected-Target Analytik
- Integration der Stoffdaten aus Stoffvollzügen (REACH, Biozid-Richtlinie...)
- Einbindung von STOFF-IDENT in die Arbeitsplattform FOR-IDENT zur Verknüpfung mit anderen Recherche-Tools (MetFrag, MassBank etc.)

Datenqualität und Datenumfang

- Integration weiterer gewässerrelevanter Stoffe in die Datenbank sowie die Aktualisierung bereits vorhandener Stoffgruppen (neu registrierte Substanzen aus Stoffvollzügen)
- Kontinuierliche Prüfung und Optimierung der Datenqualität
- Fehlersuche und -korrektur (seit Projektbeginn Umwandlung von >1.500 Fehlern und falschen Zusatzinformationen in Datensätzen)
- Konzept eines effizienten Datenmanagementsystems

Datenumfang *



Prüfregeln zur Sicherung der Datenqualität bei neuen Daten

- Substanz muss CAS-Nummer oder SMILES-Code haben
- Richtigkeit der CAS-Nummer wird über deren Prüfziffer geprüft
- SMILES und Summenformel dürfen keinen Punkt und kein * enthalten
- + im SMILES führt zu gesonderter Überprüfung der Substanzen
- $-20 < \log P < 20$

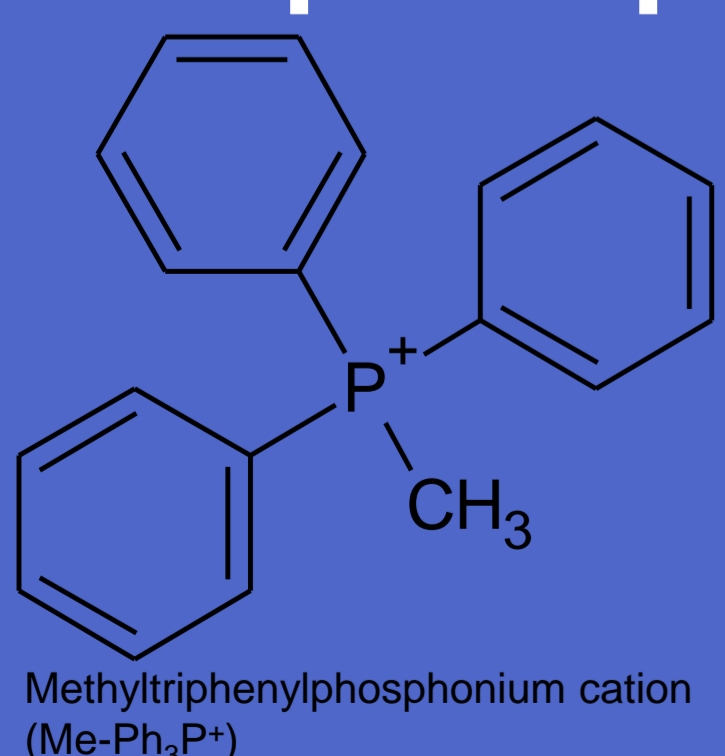
Automatisierte Fehlererkennung im
SI-Crawler, dem Daten-Einlesetool

Completen...	Suspicious	Name	Source	CAS	Source	SMILES
●	●	Acidover	UBA	99277-99-3	UBA	Ndnc(O)zncm(COCCO)Z[H]H1
●	●	ACIPROX	UBA	51037-20-0	Wikipedia	Cc1nc(C)O(=O)N1[O-]
●	●	ACITRETIN	UBA	55979-89-9	Wikipedia	CC1=CC=C(C=C1)C(=O)C=C1
●	●	ACLARUBICIN	UBA	57876-44-0	Wikipedia, chem	CC1=C(B)1O[C@@H](O)C(=O)C1
●	●	Acriflavium Chloride	UBA	8063-24-9	chemicalbook.cz	[Cl-].N1c1cc2c3cc4N(c3)nc2c1.Clc
●	●	ACTINOQUINOL	UBA	15301-40-3	http://www.drug	CCO1ccc(cc1C2ccc3C(=O)N3)O1=O
●	●	ADAPALENE	UBA	106685-40-9	Wikipedia, chem	CC1=CC=C(C=C1)C2=CC=C(C=C2)C1
●	●	Adefovirdipivoxil	UBA	142340-99-6	UBA	CC(C)C(C)C(=O)COP(=O)(O)CC(C)C
●	●	Ademetionine disulfate tosylate	UBA	97540-22-2	chemicalbook.cz	OS(=O)(=O)C(=O)O[S](=O)(=O)C
●	●	ADRENALONE	UBA	99-45-6	chemicalbook.cz	CNCC1=O[C@@H](O)C1=O
●	●	Agomelatine	UBA	138112-76-2	UBA	CCOC1=CC=C(C=C1)C=C1
●	●	AMALGAMINE	UBA	483-04-5	chemicalbook.cz	CO(C)C1=C(C)C(C)C1
●	●	AMALINE	UBA	480-12-07	chemicalbook.cz	
●	●	ALATROFLOXACIN	UBA	157605-25-9	wikipedia, chemi	CS1O(=O)=O.C1C(B)N(C1=O)N(C)C
●	●	ALBENDAZOLE	UBA	54965-21-8	Wikipedia	CCCS1ccc2nc(NC1=O)OC(=O)H2c1
●	●	ALCLOMETASONE	UBA	66734-13-2	Wikipedia	
●	●	ALCLOXA	UBA	1337-25-5	chemicalbook.cz	

Kategorisierung und Taggen

- Jede in der Datenbank enthaltene Substanz ist einer oder mehreren Kategorien zugeteilt (REACH, Pharmaceuticals, Biocides, PSM, TP's)
- Taggen bedeutet, dass eine Substanz mehreren Ursprungslisten zugewiesen werden kann, z.B. zur Spezialsuche in einzelnen Listen: REACH, einzelne Laborlisten, Liste positiv geladener Substanzen, ...

Bsp.: Triphenylphosphonium mit M suchen



Suche über M [+H]

Kein Ergebnis bei Suche im positiven Ionenmodus, da Masse um +H falsch liegt

Suche über M [±0]

Suche erfolgreich im Ionenmodus [±0], Substanz wurde gefunden, da positive Grundladung

neu Automatisierte Gruppensuche aller Substanzen mit positiver Grundladung über M [±0] möglich

Fazit und Ausblick

- Stetig steigende Nutzerzahlen (dadurch weitere Fehlerelimination)
- Verstärkter Fokus auf Transformationsprodukte
- Optimierung der automatischen Fehlererkennung
- Kontinuierliche Erhöhung des Datenbestandes durch Integration weiterer Laborlisten (national und international)

