



STOFF-IDENT- Datenbank

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

AP3 – Datenqualität und Datenumfang

- Integration weiterer gewässerrelevanter Stoffe in die Datenbank
- Aktualisierung vorhandener Stoffgruppen
- Prüfung des Datenumfangs bei der Anwendung von STOFF-IDENT
- **Kontinuierliche Prüfung und Optimierung der Datenqualität**
- **Fehlersuche und -korrektur**

Bisher in STOFF-IDENT enthaltene Stoffe

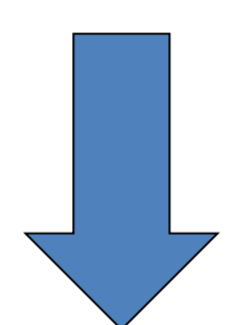
Stoff	Quelle	Anzahl	Jahr
REACH	ECHA	ca. 4350	bis Dez 2013
Arzneimittel	UBA	ca. 1500	2009
PSM & PSM-TP	Bundesamt für Verbraucherschutz und Lebensmittelsicherheit, Reemtsma et al. 2013	ca. 400	2014 in D zugelassen
Biozide	Commission Regulation (EC) Anhang I & II	ca. 300	2007
Weitere schon nachgewiesene Stoffe & TP	Zweckverband Landeswasserversorgung (LW)	2665	Feb 2014
	Fachpublikationen	ca. 1000	
	NORMAN: Candidates, List of Emerging Pollutants	ca. 870	2013
	Forensik Innsbruck (Oberacher 2011)	1207	2011
Gesamt alt (abzüglich der Überschneidungen)		ca. 8500	

Neue in STOFF-IDENT enthaltene Stoffe

Stoff	Quelle	Anzahl	Jahr
Chemikalien	BfG, Schlüsener et al.	930	2015
Arzneimittel	UBA	210	
Haushaltschemikalien	Drewes et al.	30	2009
Haushaltschemikalien	High Production Chemical Database/Colorado School of Mines/Southern Nevada Water Authority	630	
Haushaltschemikalien	Huntscha et al.	110	2012
TP's	GWA (Fenner et al. 2011)	75	2011
Gesamt neu (abzüglich der Überschneidungen)		ca. 9100	

Prüfregeln zum Einlesen in die Datenbank

- Substanz muss CAS-Nummer oder SMILES-Code haben
- Richtigkeit der CAS-Nummer wird über deren Prüfziffer überprüft
- SMILES und Summenformel dürfen keinen Punkt enthalten
- SMILES darf kein * enthalten
- $-20 < \log P < 20$



Automatisierte Fehlererkennung im **SI-Crawler**, dem Daten-Einlesetool

Completeness	Suspicious	Name	Source	CAS	Source	SMILES
●	●	Aciclovir	UBA	5077-89-3	UBA	Nc1c(O)c2nc(COCCO)c2[nH]1
●	●	ACIPIMOX	UBA	5107-30-0	Wikipedia	Cc1ccc(Oc2ccccc2)cc1
●	●	ACETHEIN	UBA	5579-83-9	Wikipedia	Cc1ccc(Cc2cc(C)cc(C)cc2)cc1
●	●	ACCLARUBIN	UBA	3728-44-0	Wikipedia, chem	CC1=C(C(=O)O)C=C(C)C=C1
●	●	Aciclovirum Chloride	UBA	8063-24-9	chemicalbook.cz	[Cl-].[Na+]ClC1=CC=C(C=C1)N2C=CC(=O)N2
●	●	ACTINOQUINOL	UBA	15301-40-3	http://www.drug	CC1=CC=C(C=C2C1=O)C=C2
●	●	ADAPALENE	UBA	10685-40-9	Wikipedia, chem	Cc1ccc(Cc2cc(C)cc(C)cc2)cc1
●	●	Adelovirginol	UBA	14230-99-6	UBA	CC1=C(C(=O)O)C=C(C)C=C1
●	●	Ademionone disulfate tosylate	UBA	9750-22-2	chemicalbook.cz	OS(=O)(=O)O[S](=O)(=O)O[S](=O)(=O)Cl
●	●	ADRENALONE	UBA	99-45-6	chemicalbook.cz	CNC1=CC=CC=C1O
●	●	Agometalin	UBA	13812-76-2	UBA	Cc1ccc(Cc2cc(C)cc(C)cc2)cc1
●	●	ANALIZINE	UBA	433-64-5	chemicalbook.cz	COC1=CC=C(C=C1)C=C(C)C=C1
●	●	ANALINE	UBA	150-13-0	chemicalbook.cz	Nc1cccnc1
●	●	ALATROLOACIN	UBA	15705-25-9	wikipedia, chemi	CS(=O)(=O)O.CC(C)N(C)C1=CC=C(C=C1)O
●	●	ALBENDAZOLE	UBA	54863-21-8	Wikipedia	CC1=C(C(=O)N)C=C(C)C=C1
●	●	ALCLOMETASONE	UBA	6073-13-2	Wikipedia	CC1=C(C(=O)N)C=C(C)C=C1
●	●	ALICLOIA	UBA	1197-25-5	chemicalbook.cz	

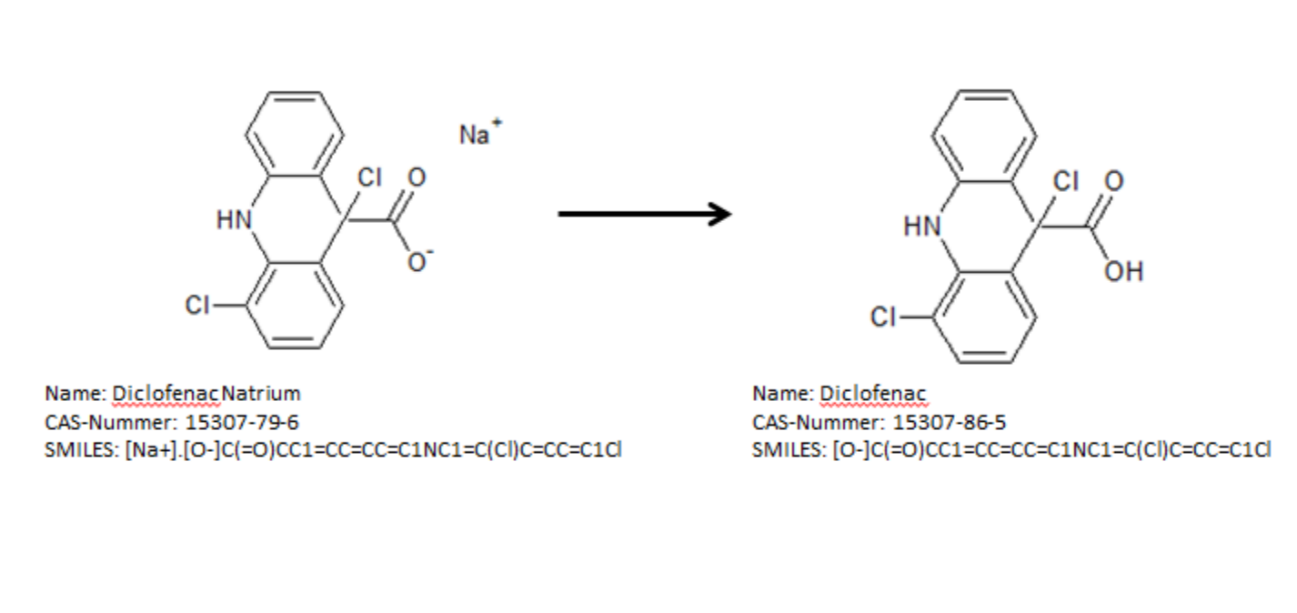
Sicherung der Datenqualität bei neuen Daten

- Manuelle Korrektur der durch die Prüfregeln angezeigten Fehler
- Nichtübereinstimmung von CAS, SMILES und IUPAC bei nochmals eingelesenen Datensätzen → Korrektur
- Manuelle Fehlerkorrekturen in neuen Datensätzen, z.B. Klärung der CAS-Nummer durch Recherchen über CAS Registry
- Nachvollziehbare Dokumentation aller Bearbeitungsschritte und Datenquellen
- Schaffung konsistenter Datensätze: Salze in Säuren umwandeln d.h. Korrektur von Masse, SMILES und CAS-Nummer

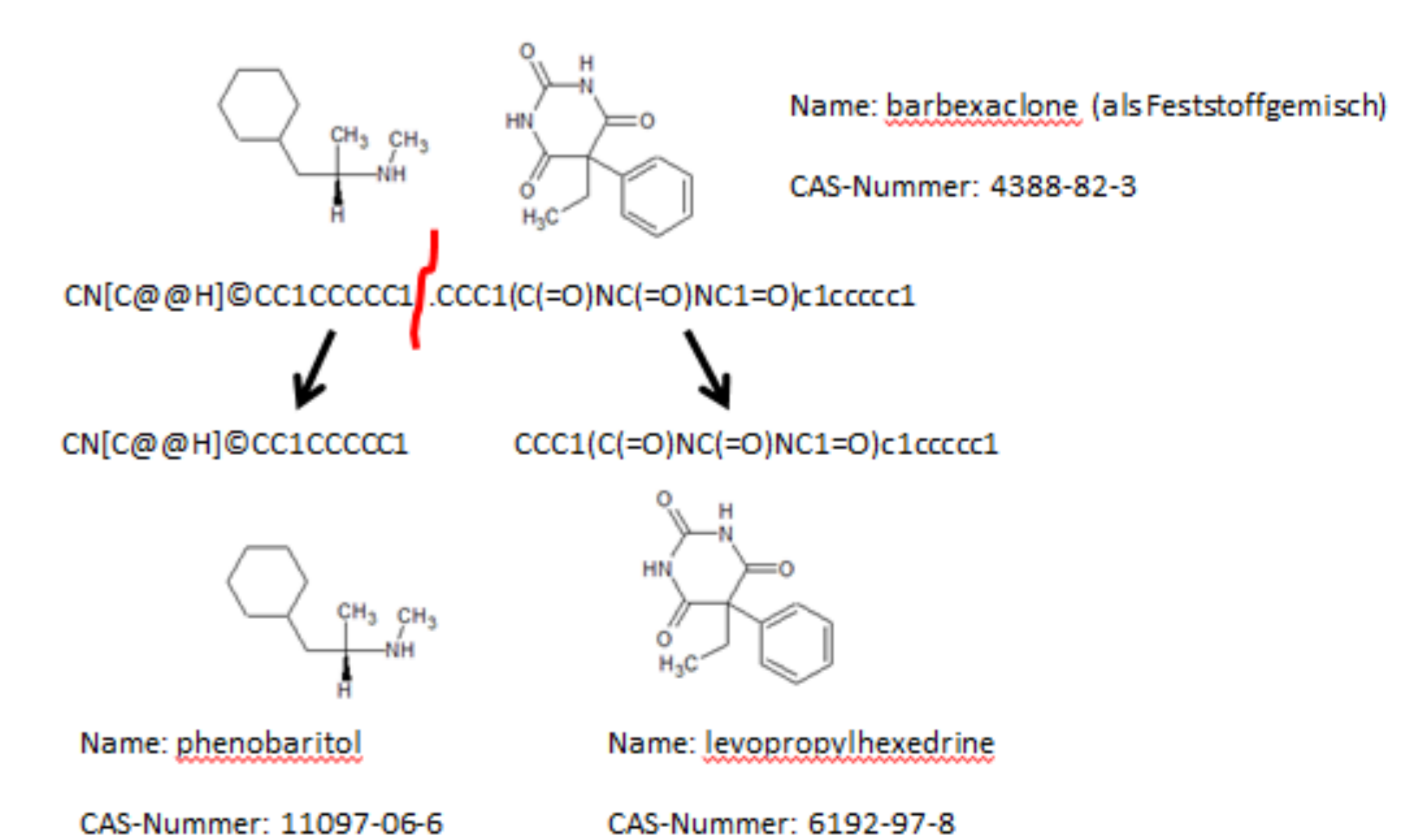
Fehlersuche

- Salze, Komplexe und Gemische
 - Punkte in über 120 SMILES-Codes
 - Punkte in der Summenformel
 - Mehrere CAS-Nummern bei über 60 Substanzen (nur eine der CAS-Nummern ist tatsächlich registriert)
 - Substanzen ohne SMILES-CODE in STOFF-IDENT
 - Falsche Übersetzung von Sonderzeichen in ?
- manuelle Korrektur von ca. 900 „Fehlern“

Beispiel: Salze



Beispiel: Punkt im SMILES



Beispiel: mehrere CAS-Nummern

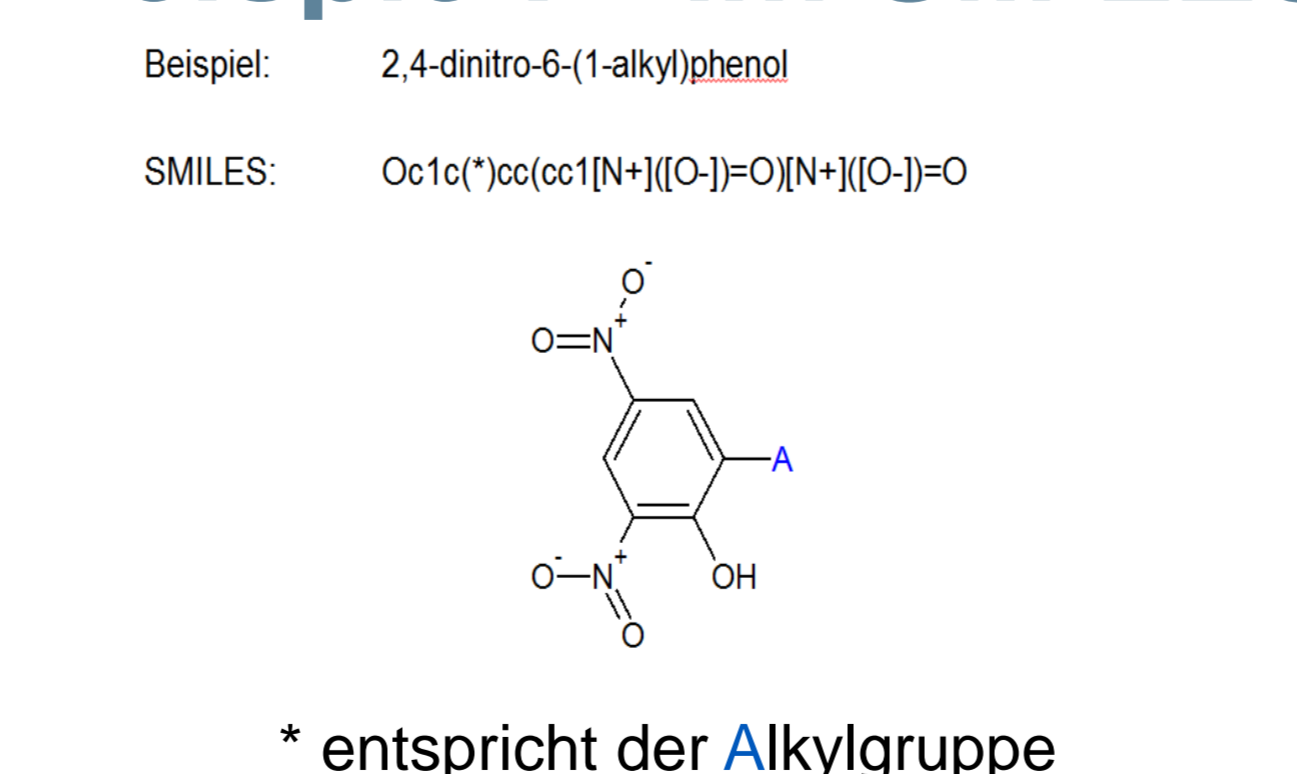
(12262-73-6/RN) registered CAS-number

```

=> d L3 RBG
1 12262-73-6, 43136-35-2, 136108-93-5, 245654-34-6, 623158-96-3,
DR 85658-15-5, 878903-72-1, 890704-54-8, 896506-46-0, 906507-37-7,
1192555-95-5
deleted CAS-numbers

```

Beispiel: * im SMILES



Ausblick

- Fokus auf Eingabe von Transformationsprodukten
- Optimierung der Kategorisierung
- Laufende Aktualisierung der Datenbank (z.B. neu unter REACH registrierte Chemikalien)

