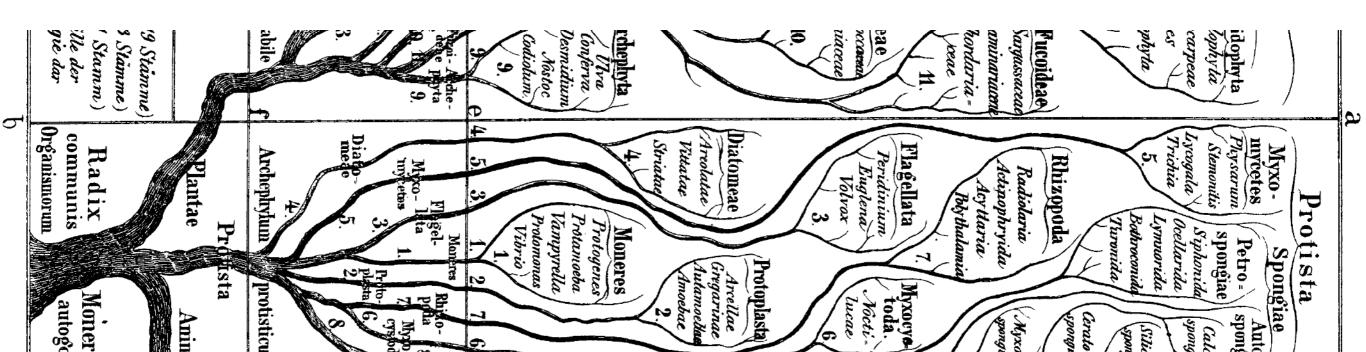


Référentiels taxinomiques et formats de données

Retour d'expérience sur EcoTaxa



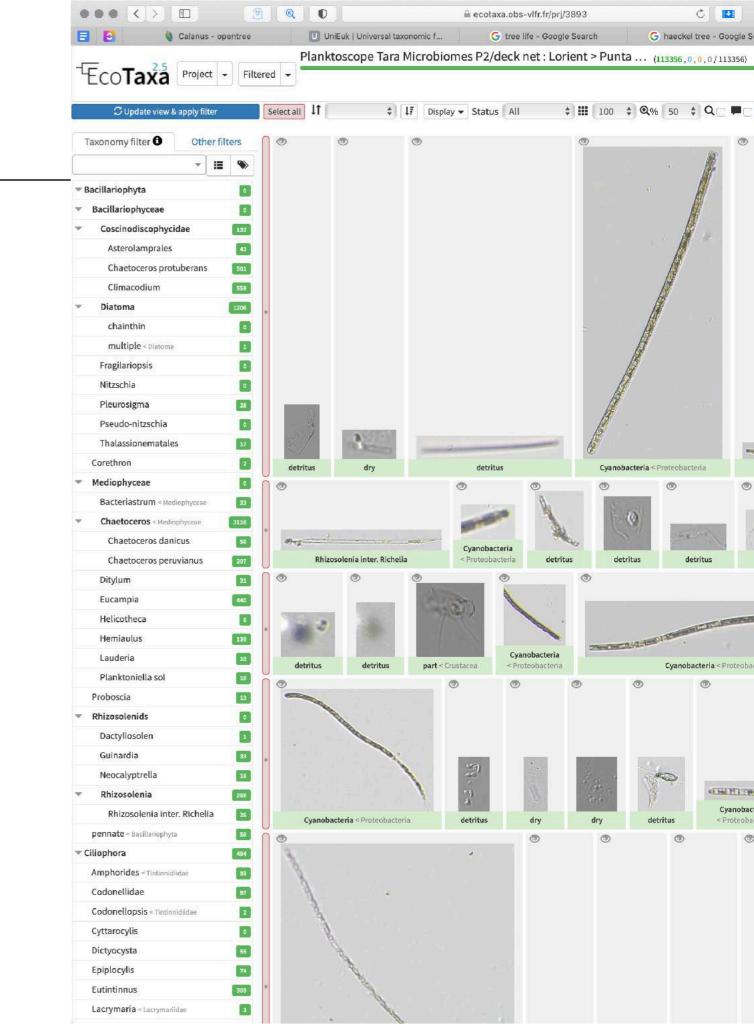
EcoTaxa

Une plateforme pour **stocker** et **classer** taxinomiquement des **images** d'organismes individuels (principalement du plancton)

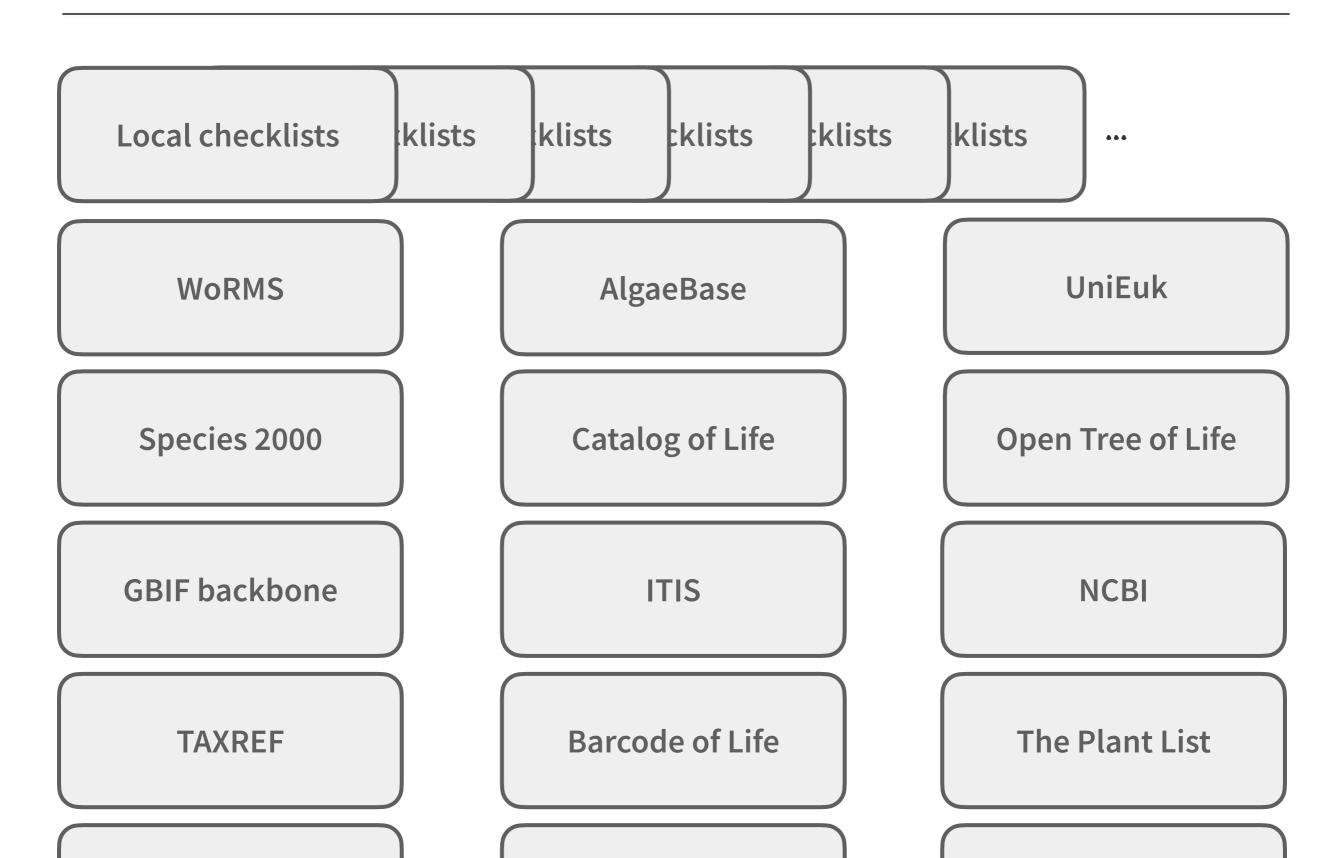
Démarré comme un outil pour les données de Tara Océans

Contient actuellement ~150M d'objets, 63M nommés, en ~45,000 points de l'océan, collectés par ~350 organisations au niveau international

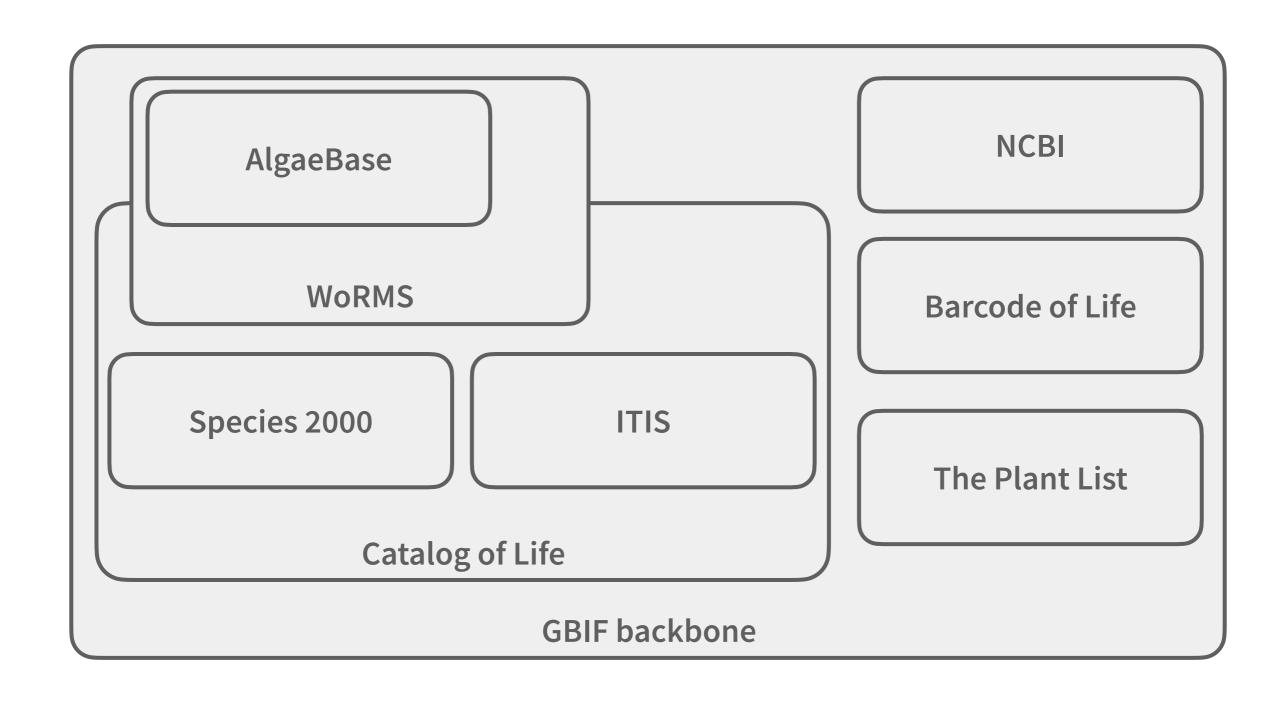
⇒ Nécessite un référentiel taxinomique universel et des formats d'import/export.



Les référentiels taxinomiques sont nombreux...



...leurs relations sont complexes



Les "référentiels" taxonomiques sont en fait toutes des méta-référentiels



Grandes divergences à la racine de l'arbre...

Catalogue of Life

▼ unranked: Biota

kingdom: Animalia

kingdom: Archaea

kingdom: Bacteria Cavalier-Smith, 2002

kingdom: Chromista

kingdom: Fungi

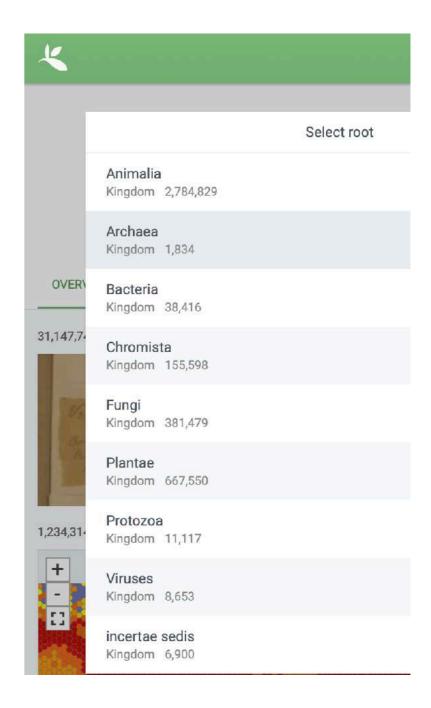
kingdom: Plantae

kingdom: Protozoa

kingdom: Viruses



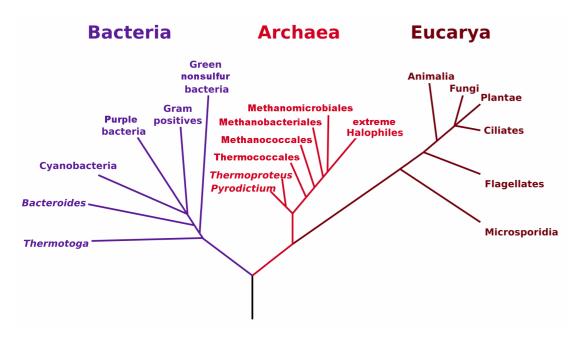
- Biota (238 252)
 - + Kingdom Animalia (203 311)
 - + Kingdom Archaea (116)
 - + Kingdom Bacteria (2 077)
 - + Kingdom Chromista (20 237)
 - + Kingdom Fungi (1 382)
 - + Kingdom Plantae (10 378)
 - + Kingdom Protozoa (633)
 - + Kingdom Viruses (115)

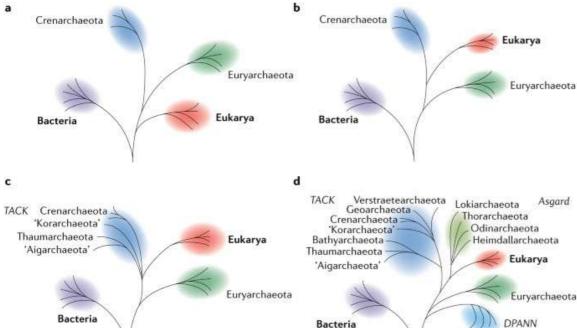


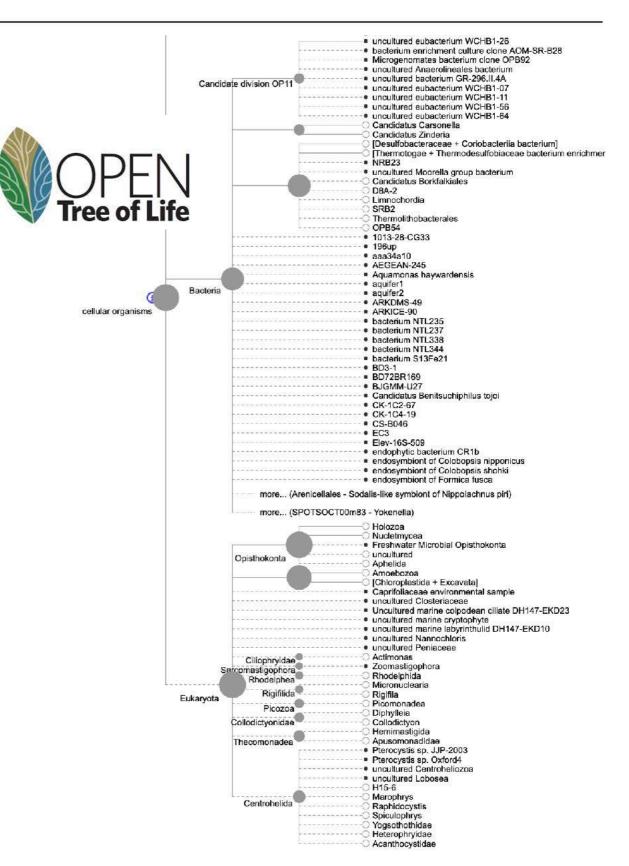
Grandes divergences à la racine de l'arbre...

Nature Reviews | Microbiology

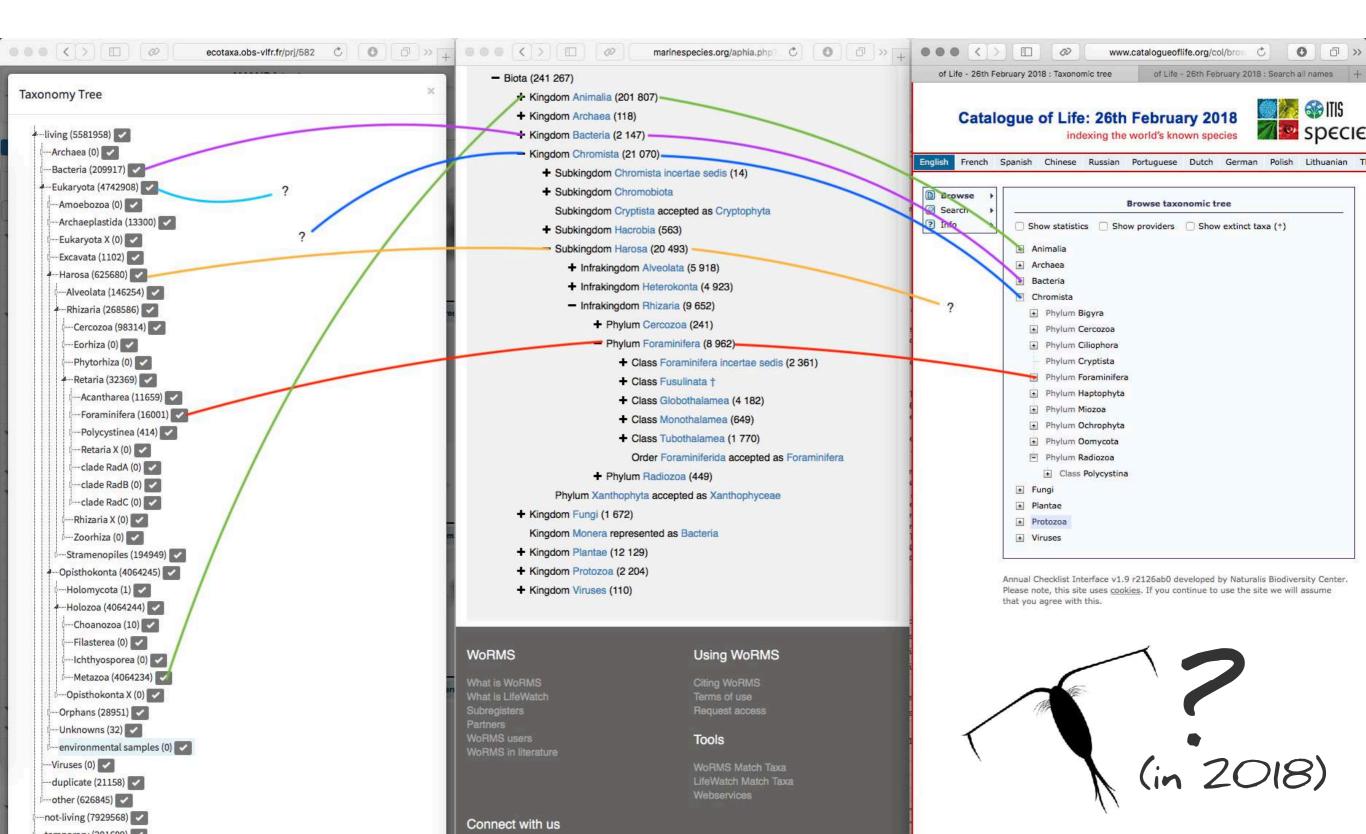








... et sur la topologie de l'arbre



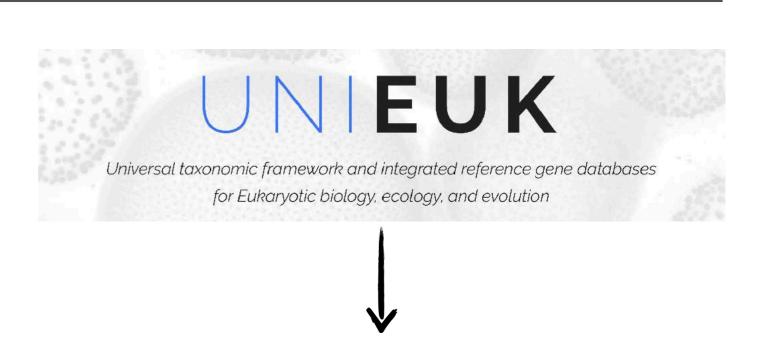
Et dans EcoTaxa, finalement?

Taxonomie initiale par **UniEuk**, pour des raisons historiques

Ajouts, par les utilisateurs, de taxa manquants (et de "taxa" morphologiques)

Transition vers **WoRMS**, pour compatibilité avec OBIS

À terme, espérons l'émergence d'un **méta**réferentiel taxinomique assurant la correspondance avec l'existant, universel, phylogénétiquement correct...



World Register of Marine Species

Le format d'EcoTaxa est peu contraignant

Dossier d'images + 1 fichier .tsv avec une ligne par image (donc bcp de répétitions)

Organisation hiérarchique des métadonnées

Seuls deux champs sont obligatoires

Certains champs ont un **format** imposé

La **plupart** sont **libres** et juste stockés par dataset

⇒ Flexible pour les utilisateurs mais peu standard

IMAGE

- o img_file_name [t]: name of the image file in the folder (including extension)
- img_rank [f]: rank of image to be displayed, in case of existence of multiple (<)
 Starts at 1.</pre>
- OBJECT: one object to be classified, usually one organism. One object can be represented tsv file, there is one line per image which means the object data gets repeated on seven
 - o object_id [t]: identifier of the object, must be unique in the project. It will be
 - o object_link [f]:URL of an associated website
 - o object_lat [f]:latitude, decimal degrees
 - o object_lon [f]:longitude, decimal degrees
 - object_date [f]:ISO8601 YYYYMMJJ UTC
 - o object_time [f]:ISO8601 HHMMSS UTC
 - object_depth_min [f]: minimum depth of object, meters
 - object_depth_max [f]: maximum depth of object, meters

And, for already classified objects

- object_annotation_date [t]:ISO8601 YYYYMMJJ UTC
- object_annotation_time [t]:ISO8601 YYYYMMJJ UTC
- object_annotation_category [t]: class of the object with optionally its dir by left angle bracket without whitespace "Cnidaria<Hydrozoa" or old style between (Hydrozoa)"
- object_annotation_category_id [f]: Ecotaxa ID of the class of the object, export
- object_annotation_person_name [t]: name of the person who identified t
 object_annotation_person_email [t]: email of the person who identified
- object_annotation_status [t]: predicted, dubious, or validated

And additional object-related fields

- object_*** [f] or [t]: other fields relative to the object. Up to 500 [f] fields and
- PROCESS: metadata relative to the processing of the raw images
 - process_id [t]: identifier. The processing information is associated with the a missing, a dummy processing identifier will be created.
 - o process_*** [t]: other fields relative to the process. Up to 30 of them.
- ACQUISITION: metadata relative to the image acquisition
 - acq_id [t]: identifier of the image acquisition, must be unique in the project. If identifier will be created.
 - o acq_instrument [t]: name of the instrument (UVP, ZOOSCAN, FLOWCAM, etc.
 - o acq_*** [t]: other fields relative to the acquisition. Up to 30 of them.
- SAMPLE: a collection event
 - sample_id [t]: identifier of the sample, must be unique in the project. If missir will be created.
 - o sample_*** [t]: other fields relative to the sample. Up to 30 of them.

Ajout d'un export au format DwCA

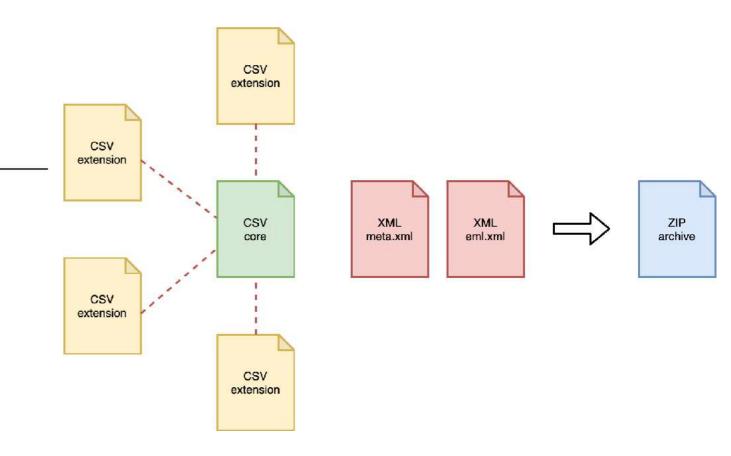
Format beaucoup plus compliqué

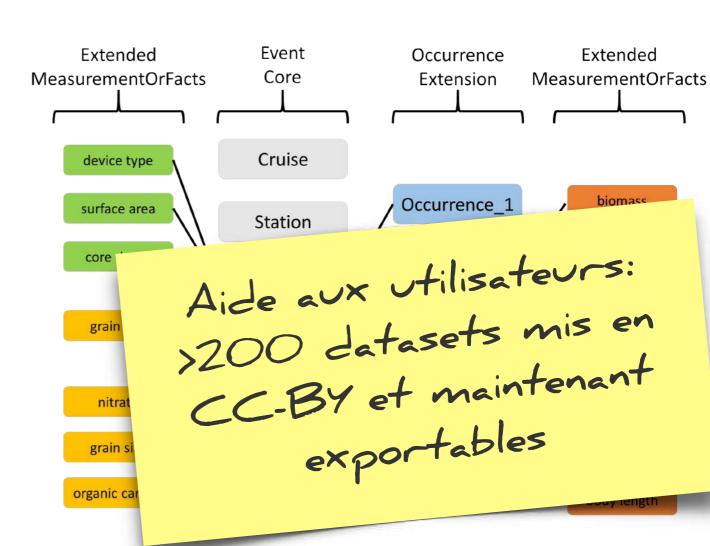
Format beaucoup plus verbeux

Format beaucoup plus **standardisé** (termes DarwinCore, vocabulaire du BODC, etc.)

⇒ Adapté pour l'échange de machine à machine mais peu agréable pour travailler







Considération du format BioODV

Format "**tableur**" relativement classique = lisible par ODV mais aussi "Excel" (et R, et Python, etc.)

Format relativement **compact** (car les champs répétés peuvent être omis)

Format pouvant être **standardisé** (par des références au vocabulaire BODC dans l'**en-tête**)

⇒ **Bon compromis** entre flexibilité et standardisation

Δ	А	В	C	D	E	F	G	Н	ı	J	K	L	IVI	N	U
1	//														
2	// <sdn_refer< td=""><td>ence xlink:h</td><td>ref="http://se</td><td>adata.bsh.de/cg</td><td>i-csr/XML/x</td><td>xmlDownload</td><td>_V2.pl?edmo</td><td>=269&identif</td><td>er=GN36199</td><td>704604" xlink</td><td>c:role="isObse</td><td>rvedBy" xlin</td><td>k:type="SDN:L</td><td>23::CSR"/></td><td></td></sdn_refer<>	ence xlink:h	ref="http://se	adata.bsh.de/cg	i-csr/XML/x	xmlDownload	_V2.pl?edmo	=269&identif	er=GN36199	704604" xlink	c:role="isObse	rvedBy" xlin	k:type="SDN:L	23::CSR"/>	
3	// <sdn_refer< td=""><td>ence xlink:h</td><td>ref="http://vo</td><td>cab.nerc.ac.uk/c</td><td>collection/C</td><td>17/current/36</td><td>AE" xlink:role</td><td>e="isObserved</td><td>By" xlink:typ</td><td>e="SDN:L23::</td><td>NVS2CON"/></td><td></td><td></td><td></td><td></td></sdn_refer<>	ence xlink:h	ref="http://vo	cab.nerc.ac.uk/c	collection/C	17/current/36	AE" xlink:role	e="isObserved	By" xlink:typ	e="SDN:L23::	NVS2CON"/>				
4				adatanet.maris2											
5				adatanet.maris2											
6	// <sdn_refer< td=""><td>ence xlink:h</td><td>ref="http://se</td><td>adatanet.maris2</td><td>2.nl/v_cdi_v</td><td>/3/print_xml.a</td><td>sp?edmo=26</td><td>9&identifier=</td><td>GN36199764</td><td>604001_0MS</td><td>B6_269_G04'</td><td>xlink:role="</td><td>isDescribedBy</td><td>' xlink:type=</td><td>"SDN:L23::CDI"</td></sdn_refer<>	ence xlink:h	ref="http://se	adatanet.maris2	2.nl/v_cdi_v	/3/print_xml.a	sp?edmo=26	9&identifier=	GN36199764	604001_0MS	B6_269_G04'	xlink:role="	isDescribedBy	' xlink:type=	"SDN:L23::CDI"
7	//SDN_paran														
8				ject> <object>SI</object>											
9			•	ject> <object>SI</object>					• • • • • • • • • • • • • • • • • • • •						
10				subject> <object< td=""><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></object<>											
11				ubject> <object></object>									~		
12	// <subject>S</subject>	DN:LOCAL:V	anadium V <td>ubject><object></object></td> <td>SDN:P01::R</td> <td>RWSSED02<td>bject><units></units></td><td>SDN:P06::UN</td><td>1KG</td><td></td><td></td><td></td><td></td><td>(h)</td><td><u> </u></td></td>	ubject> <object></object>	SDN:P01::R	RWSSED02 <td>bject><units></units></td> <td>SDN:P06::UN</td> <td>1KG</td> <td></td> <td></td> <td></td> <td></td> <td>(h)</td> <td><u> </u></td>	bject> <units></units>	SDN:P06::UN	1KG					(h)	<u> </u>
13	//		Za cons						NAME OF TAXABLE PARTY.						
		Station	Туре	YYYY-MM-DI L	ongitude [d	Latitude [deg	LOCAL_CDI_I	EDMO_code	Bot. Depth [COREDIST [r	QV:SEADATA	Lithium Li			c QV:SEADATA
15	MTPII-MATE	MNB4	*	1997-09-01T	25.511	39.778	GN36199764	269	95		-				0 1
16										0.105	-				0 1
17			- Company							0.205			0		0 1
18	MTPII-MATE	MSB5	*	1997-09-01T	24.683	36.193	GN36199764	269	1255	1	-		0 1		0 1
19										0.105	1		0 1		0 1
20										0.205	1		0 1		0 1
21										0.26	1		0 1		0 1
22			To a constant							0.315	-		0 1		0 1
23	MTPII-MATE	MSB6	*	1997-09-01T	25.705	36	GN36199764	269	1300	0.005	1) 1		0 1













