

Importance de la sémantique pour la synthèse des données en écologie

Eric Garnier



CENTRE D'ÉCOLOGIE
FONCTIONNELLE
& ÉVOLUTIVE

UMR 5175, Montpellier

Data-intensive Science: A New Paradigm for Biodiversity Studies

STEVE KELLING, WESLEY M. HOCHACHKA, DANIEL FINK, MIREK RIEDEWALD, RICH CARUANA, GRANT BALLARD, GILES HOOKER

Données et écologie

ZooKeys 150: 15–51 (2011)
doi: 10.3897/zookeys.150.1766
www.zookeys.org

REVIEW ARTICLE

A peer-reviewed open-access journal
ZooKeys
Launched to accelerate biodiversity research

GUEST EDITORIAL GUEST EDITORIAL GUEST EDITOR

Ecological data in the Information Age

Data issues in the life sciences

Anne E. Thessen, David J. Patterson

Biodiversity data should be published, cited, and peer reviewed

Mark J. Costello¹, William K. Michener², Mark Gahegan³, Zhi-Qiang Zhang⁴, and Philip E. Bourne⁵

CONCEPTS AND QUESTIONS

Big data and the future of ecology

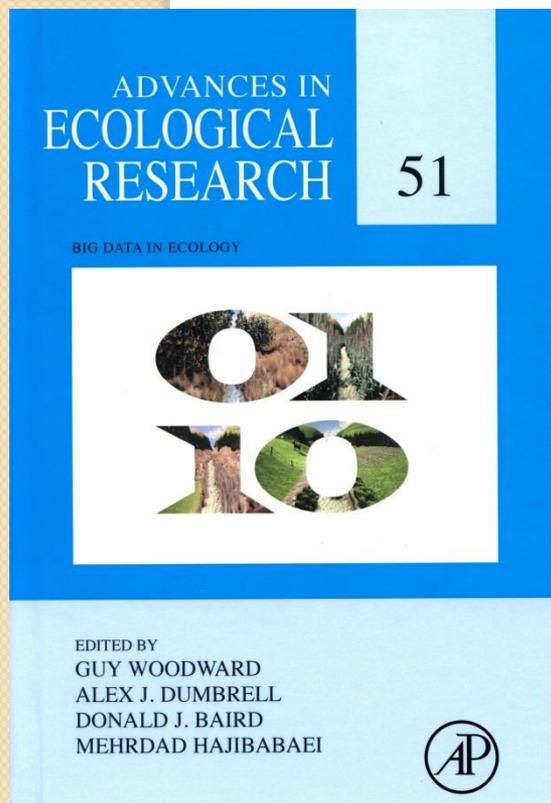
Stephanie E Hampton^{1*}, Carly A Strasser², Joshua J Tewksbury³, Wendy K Gram⁴, Amber E Budden⁵, Archer L Batcheller⁶, Clifford S Duke⁷, and John H Porter⁸

Journal of Vegetation Science 27 (2016) 865–867

COMMENTARY

Vegetation science in the age of big data

Scott L. Collins



Science

11 February 2010 | \$10



EDITORIAL

Making Data Maximally Available

INTRODUCTION

Challenges and Opportunities

PERSPECTIVE

Challenges and Opportunities of Open Data in Ecology

O. J. Reichman,* Matthew B. Jones, Mark P. Schildhauer

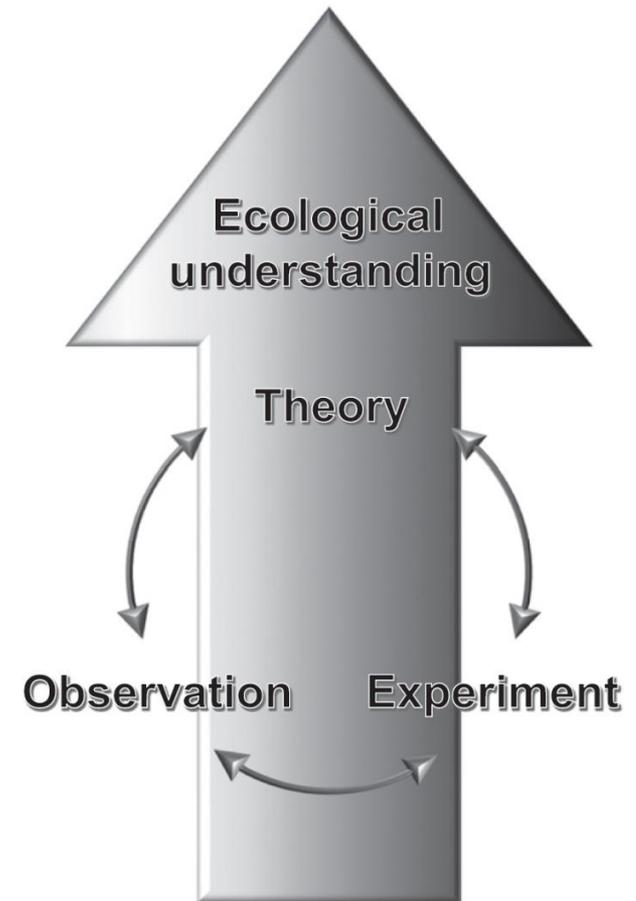
Donnée: ce qui est connu et admis, et qui sert de base, à un raisonnement, à un examen ou à une recherche

Organisation de l'exposé

- Pourquoi cet intérêt?
- Caractéristiques des données en écologie: les « données obscures »
- Exemples d'obstacles terminologiques à la synthèse des données

Pourquoi cet intérêt?

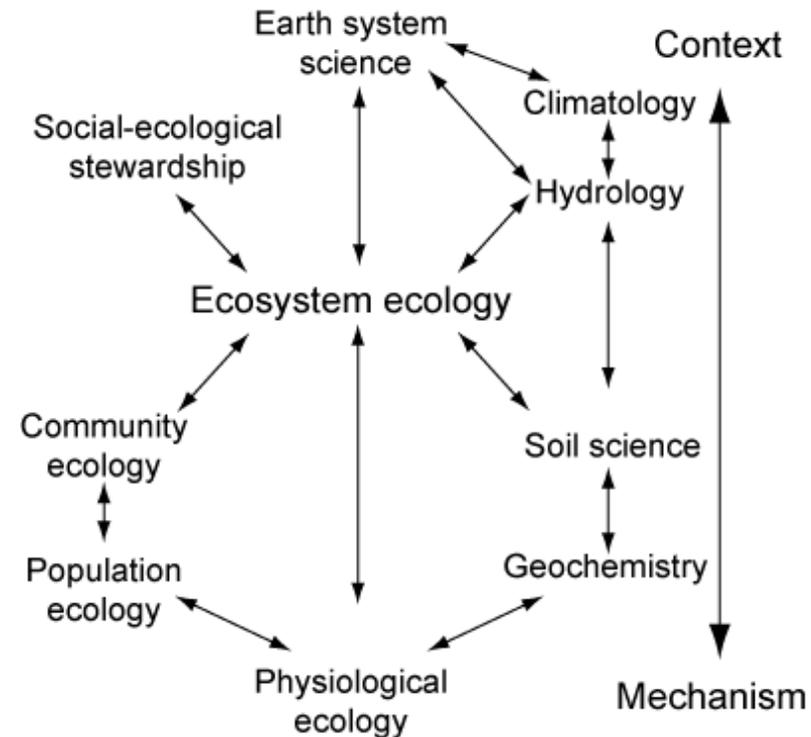
- Les données sont l'un des fondements des avancées scientifiques
- L'écologie est une discipline intégrative
- Explosion de données dans certains champs de l'écologie (« big data »)
- Certaines des questions posées à l'écologie ont une dimension continentale et/ou planétaire
- Développement des activités de synthèse scientifique



Eisenhauer *et al.*
(2016) *JVS* 27 :
1061

Pourquoi cet intérêt?

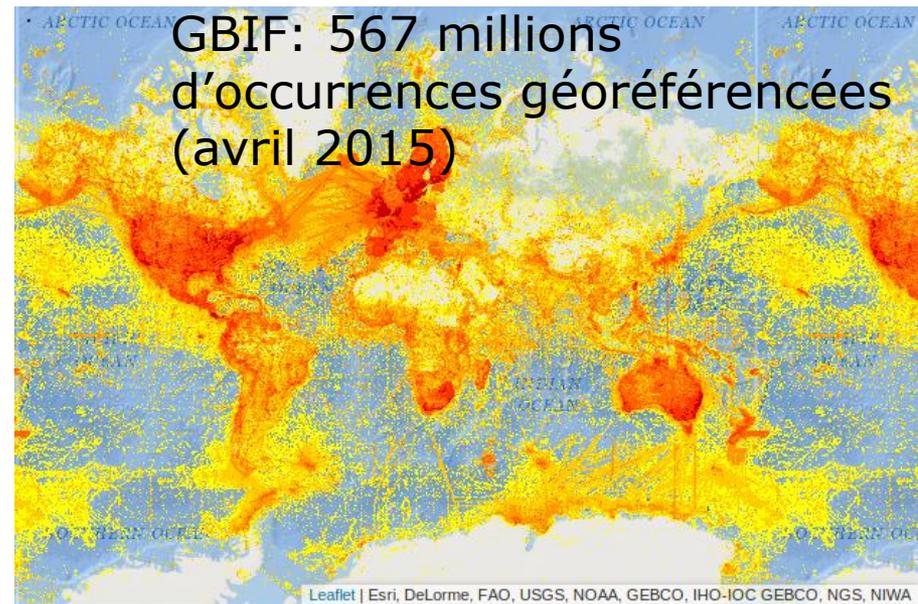
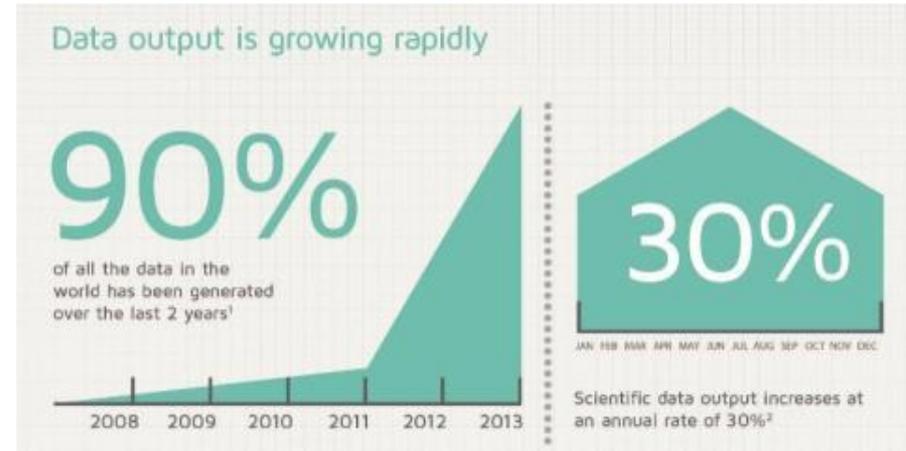
- Les données sont l'un des fondements des avancées scientifiques
- L'écologie est une discipline intégrative
- Explosion de données dans certains champs de l'écologie (« big data »)
- Certaines des questions posées à l'écologie ont une dimension continentale et/ou planétaire
- Développement des activités de synthèse scientifique



Chapin *et al.* (2011)
*Principles of Terrestrial
Ecosystem Ecology*

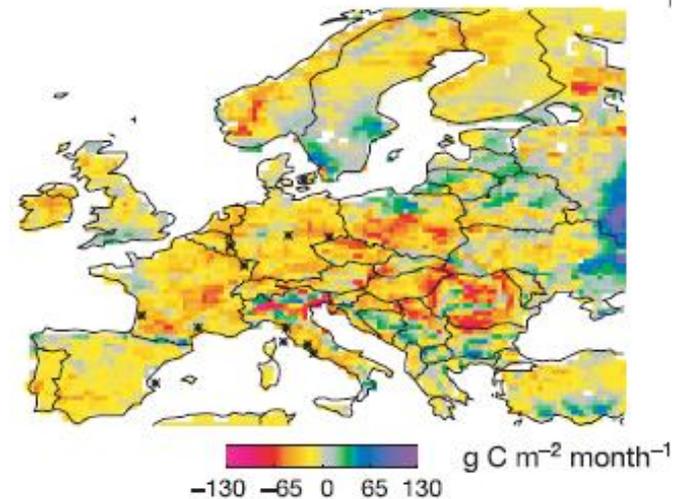
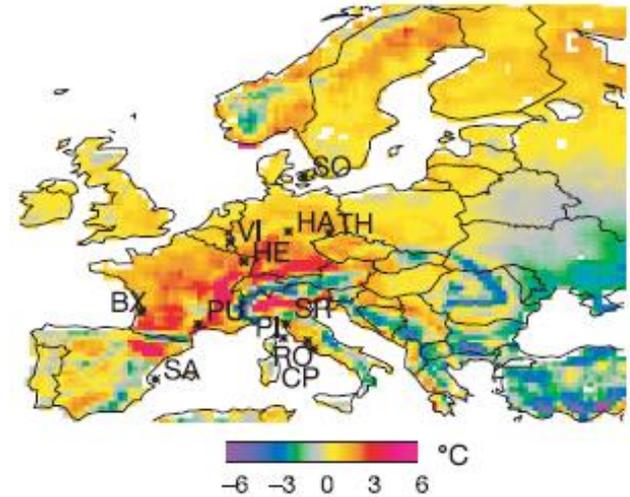
Pourquoi cet intérêt?

- Les données sont l'un des fondements des avancées scientifiques
- L'écologie est une discipline intégrative
- Explosion de données dans certains champs de l'écologie (« big data »)
- Certaines des questions posées à l'écologie ont une dimension continentale et/ou planétaire
- Développement des activités de synthèse scientifique



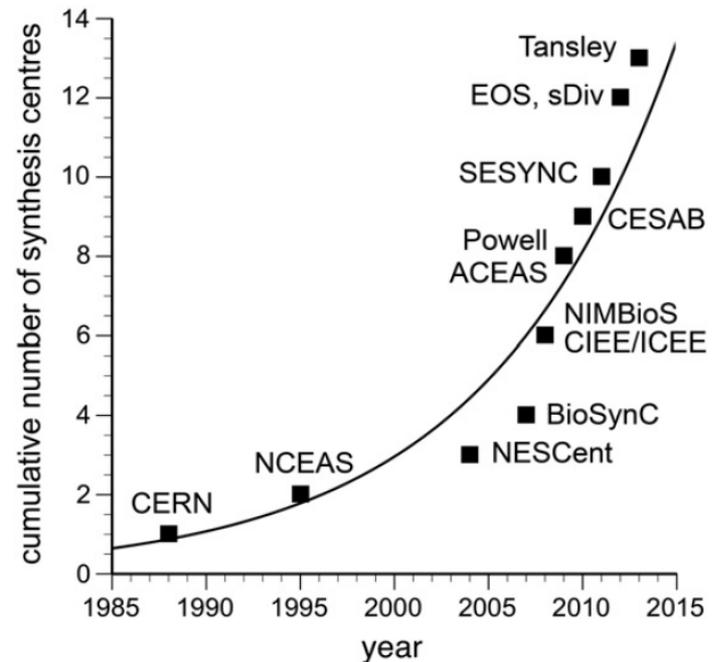
Pourquoi cet intérêt?

- Les données sont l'un des fondements des avancées scientifiques
- L'écologie est une discipline intégrative
- Explosion de données dans certains champs de l'écologie (« big data »)
- Certaines des questions posées à l'écologie ont une dimension continentale et/ou planétaire
- Développement des activités de synthèse scientifique



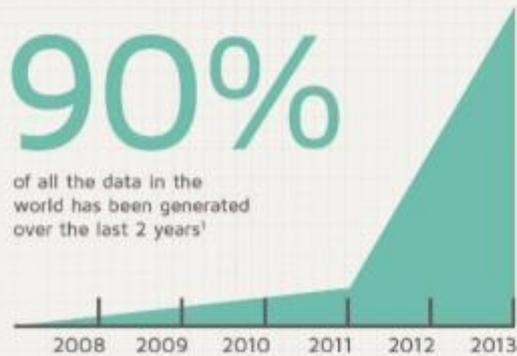
Pourquoi cet intérêt?

- Les données sont l'un des fondements des avancées scientifiques
- L'écologie est une discipline intégrative
- Explosion de données dans certains champs de l'écologie (« big data »)
- Certaines des questions posées à l'écologie ont une dimension continentale et/ou planétaire
- Développement des activités de synthèse scientifique



Des données, mais...

Data output is growing rapidly



Scientific data output increases at an annual rate of 30%²

Despite significant investment, data is not being managed effectively

\$1.5 TRILLION

is the current estimated total global spend on R&D, which could be at risk³



In one study, the odds of sourcing datasets declined by 17% each year, with 80% of datasets over 20 years old not available⁴

Funders now require data management and sharing policies



Key funding bodies such as the NIH, MRC and Wellcome Trust now request data management plans be part of applications^{5,10}

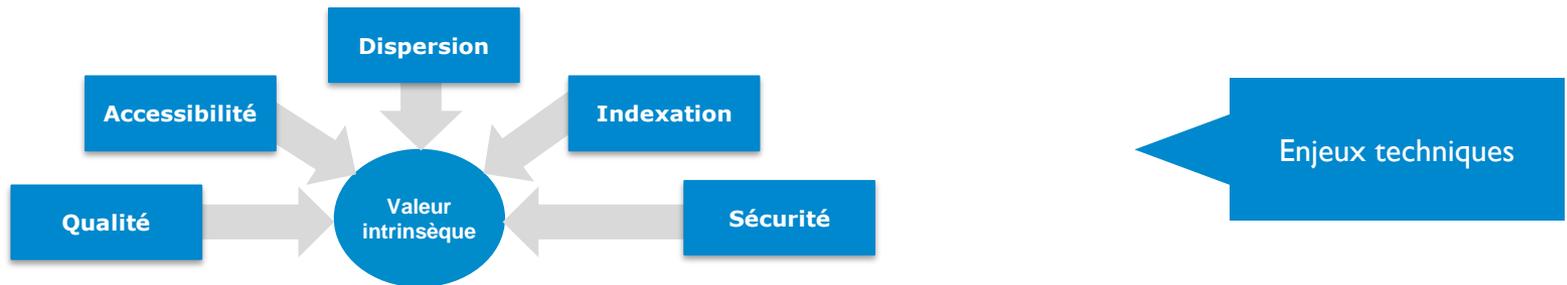
Much of the data remains unverifiable



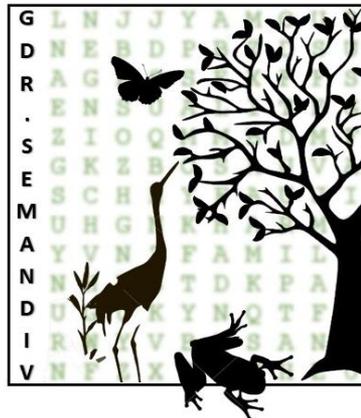
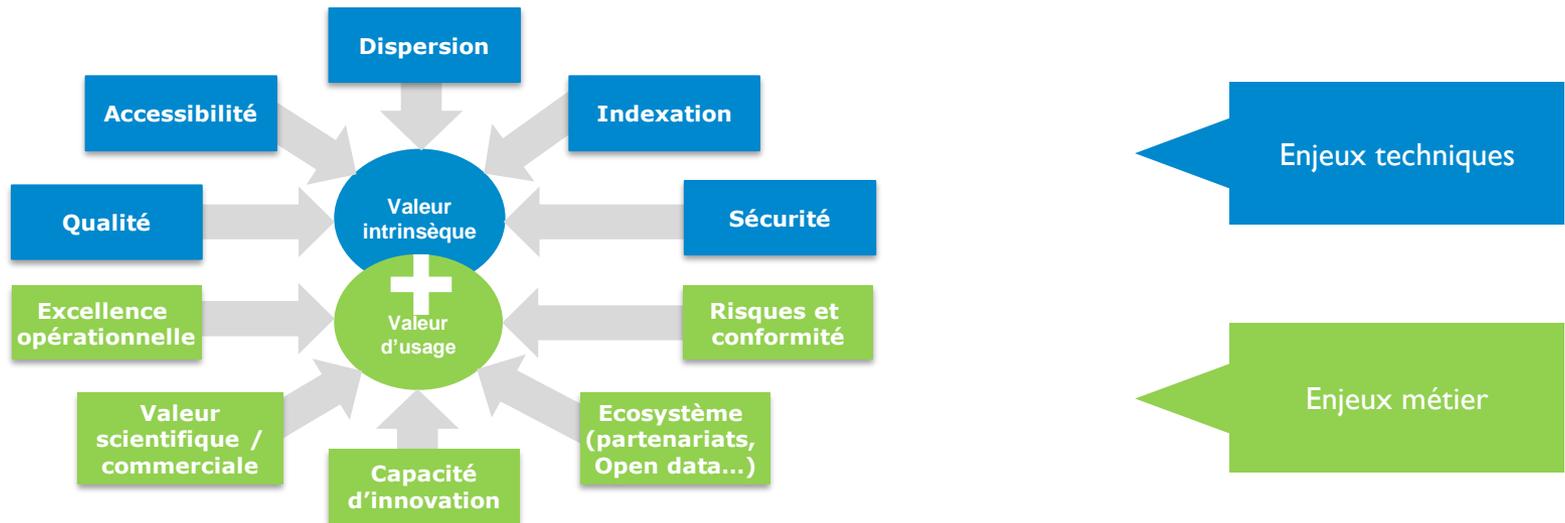
54%

of the resources used across 238 published studies could not be identified, making verification impossible⁶

Deux perspectives sur les données



Deux perspectives sur les données





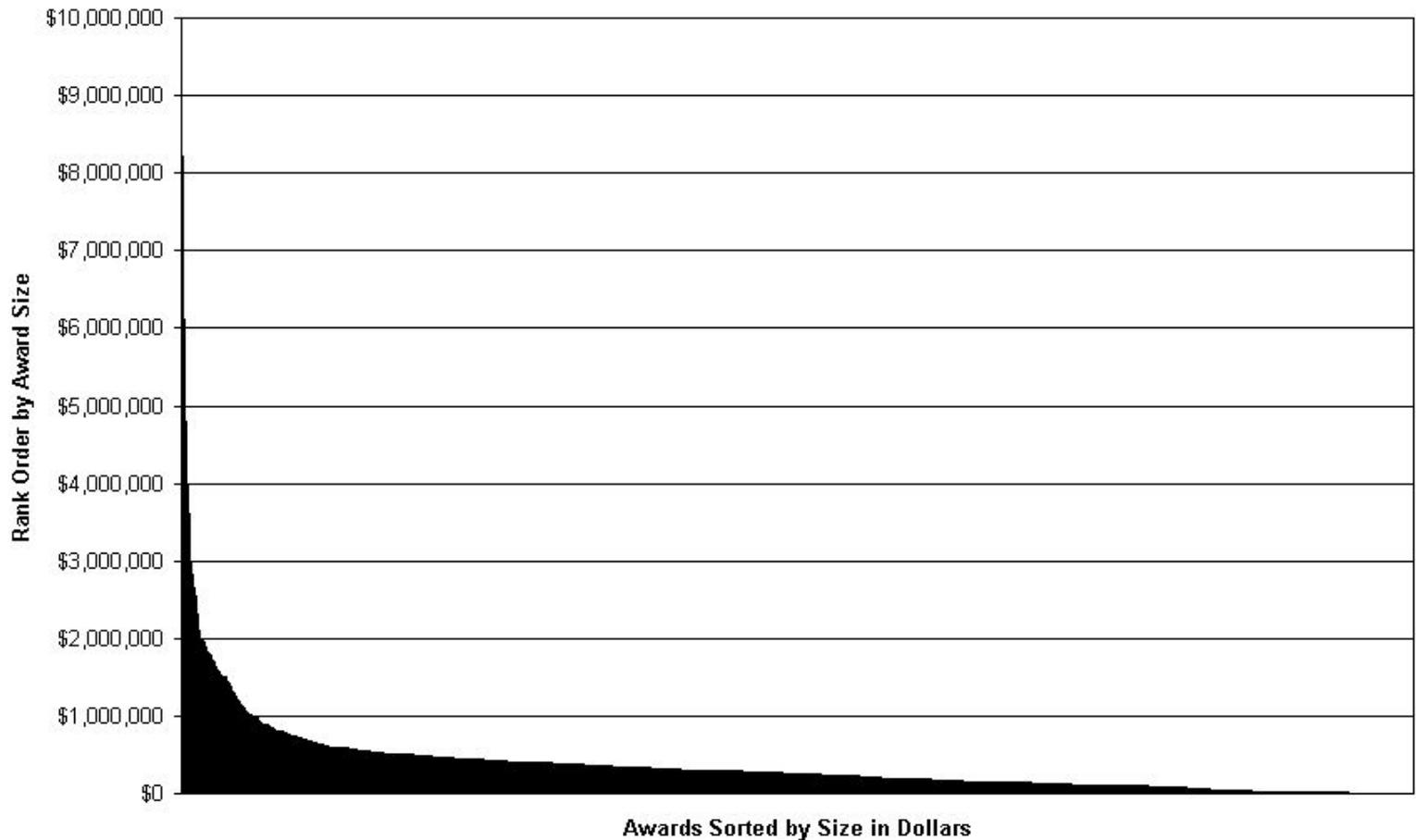
CARACTÉRISTIQUES DES DONNÉES EN ECOLOGIE

Les écologues génèrent une quantité considérable de données mais...

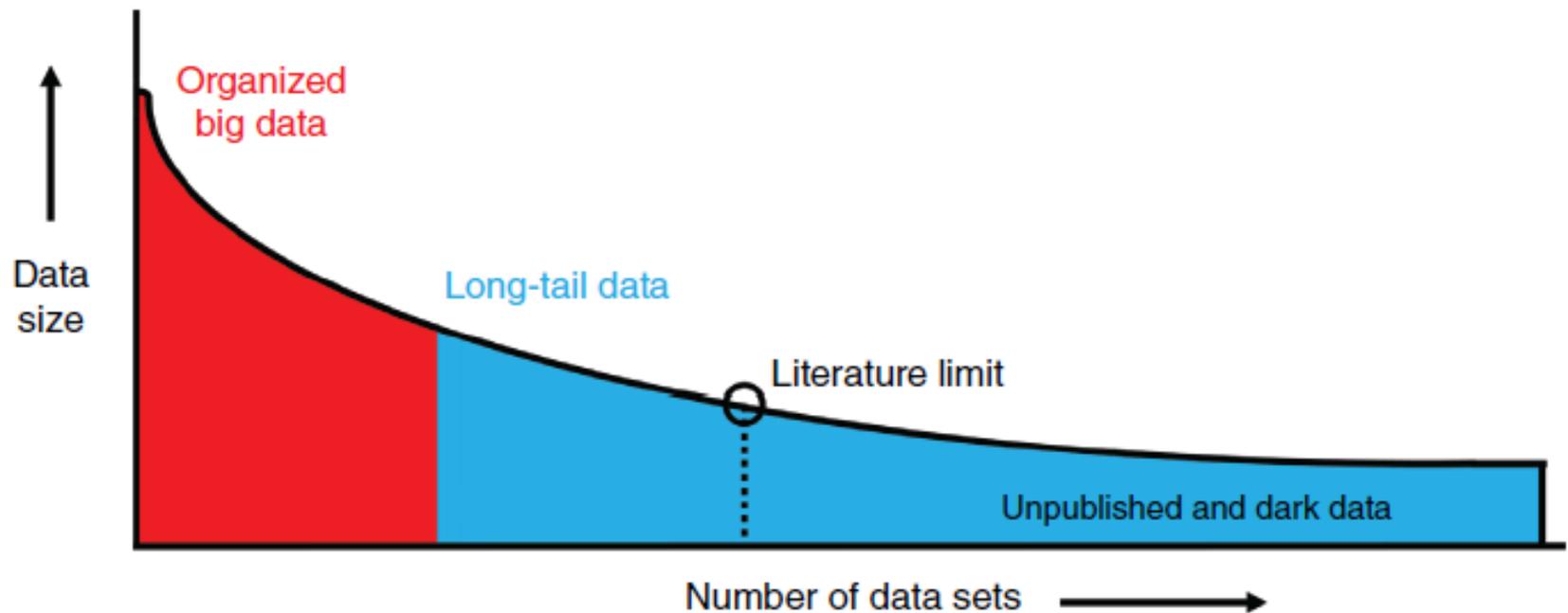
- majoritairement conduite par des individus et/ou groupes, sans réelle concertation
- peu de culture d'échange
- pas de procédures universellement admises et appliquées pour combiner et synthétiser ces données
- dominée par « long tail science » : beaucoup de petits jeux de données dispersés
 - hétérogénéité syntaxique
 - hétérogénéité sémantique

La structure du financement de la recherche...

National Science Foundation 2007 Awards



... et les « données obscures » de la « long tail science »



Ferguson *et al.* (2014) *Nature Neuro* **17** : 1142

“Donnée obscure”: toute donnée qui n’est pas facilement découverte par un utilisateur potentiel (Heidorn, 2008)

Nature des jeux de données en écologie

Caractéristiques	Type de recherche	
	Coordonnée	Indépendante
Format	Homogène	Hétérogène
Collecte et saisie	Automatisées	Manuelles
Procédures	Uniformes	Contingentes
Archivage	Centralisé	Individuel
Dépôt	Bases disciplinaires et de référence	Equipes ou laboratoires
Entretien	Réalisé	Non réalisé
Accès	Ouvert	Obscur ou non protégé
Réutilisation	Immédiate	Rare
Valorisation	Dans le cadre de l'activité professionnelle	Souvent inaperçue

Le gâchis des « données obscures »

- moins de 1% des données acquises en écologie sont accessibles après publication des résultats associés (Reichman *et al.* 2011) => opportunités perdues
- répétition inutile: coût humain et financier
- combinaison avec des données complémentaires pour répondre à des questions différentes et/ou plus larges impossible

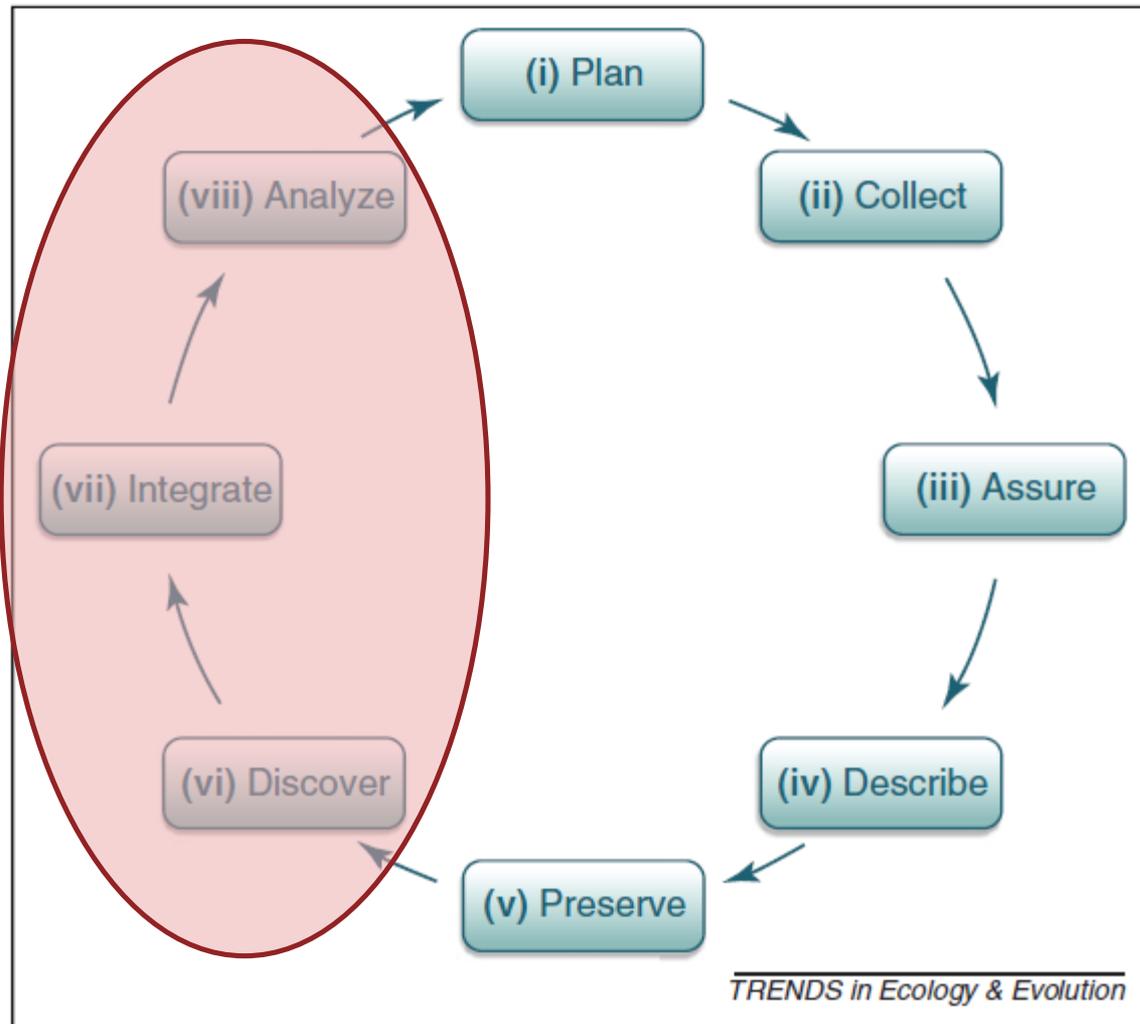


INTEGRATION DES DONNEES POUR LA SYNTHESE SCIENTIFIQUE

La synthèse scientifique

- Combinaison et intégration de différentes recherches afin d'améliorer la généralité et l'applicabilité des résultats de la recherche scientifique
- Elle peut se faire au sein d'une discipline, entre disciplines et entre secteurs d'activité professionnelle: elle n'est **pas** synonyme d'interdisciplinarité

Synthèse et cycle de vie des données

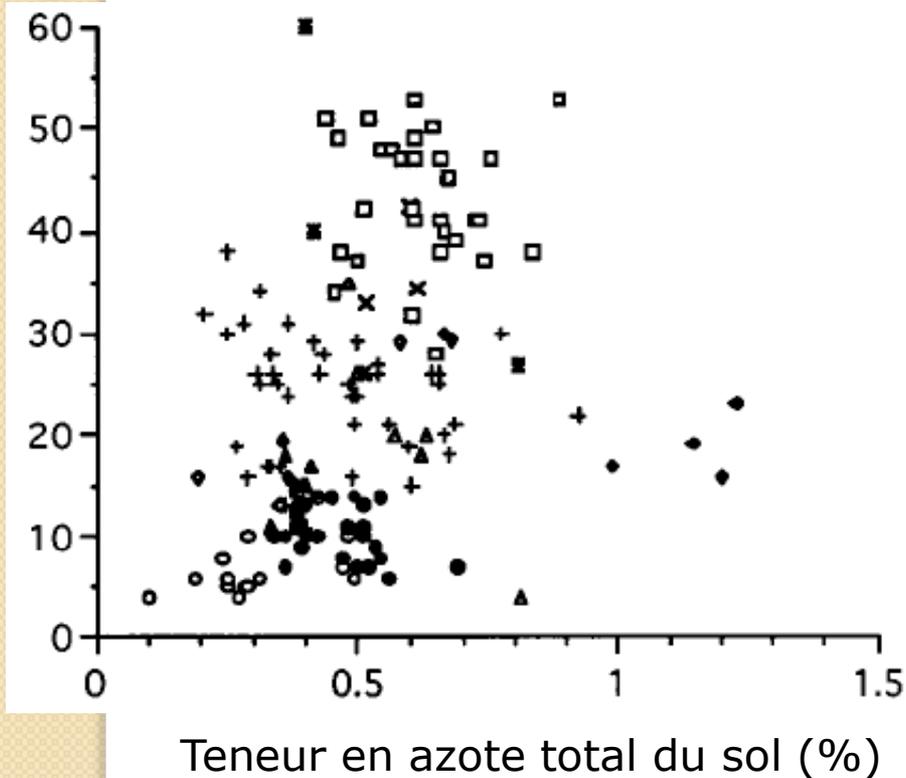


La synthèse scientifique en sciences de la biodiversité

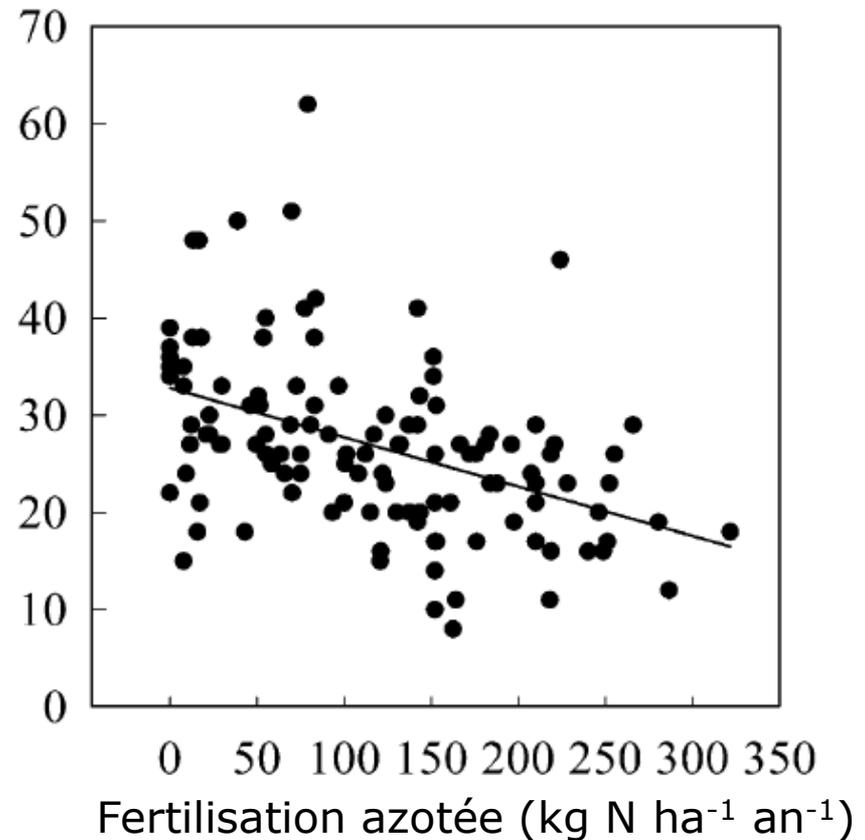
- **Les questions relatives à la biodiversité sont généralement très complexes:**
 - Concepts pas toujours bien définis
 - La robustesse statistique et la généralité des résultats de la recherche en biodiversité et écologie doivent être améliorées
- **Le problème des petits jeux de données et/ou des données obscures**
 - Découverte et accessibilité
 - Hétérogénéité sémantique et syntaxique

Ex. I: quelle est la relation entre diversité des plantes et fertilité des prairies permanentes ?

Nombre d'espèces de plantes



Janssens *et al.* (1998)
Plant Soil 202: 69



Klimek *et al.* (2007)
Biol Cons 134: 559

Les éléments nécessaires à la compréhension de cette relation

- **Définition des concepts:**

- biodiversité
- fertilité

- **Choix des variables descriptives:**

- nombre d'espèces (indicateur de biodiversité)
- composantes de la disponibilité en azote (indicateur de fertilité)

- **Méthodes**

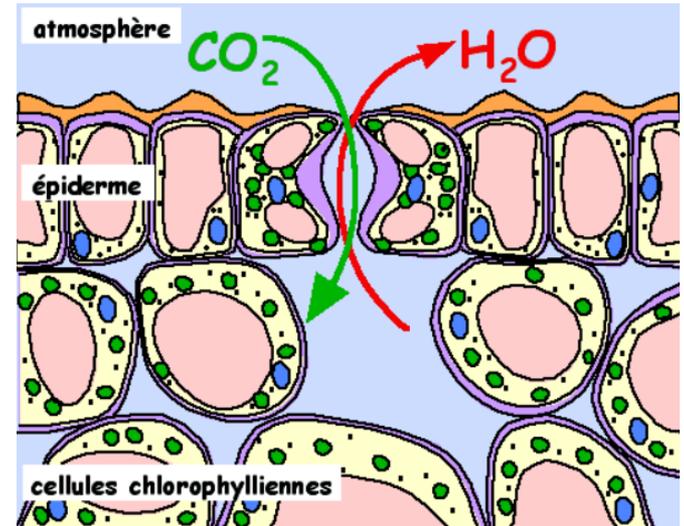
- **Interactions avec les autres facteurs:**

- contexte biogéographique
- contexte pédo-climatique
- autres ressources
- autres pratiques de gestion
- ...

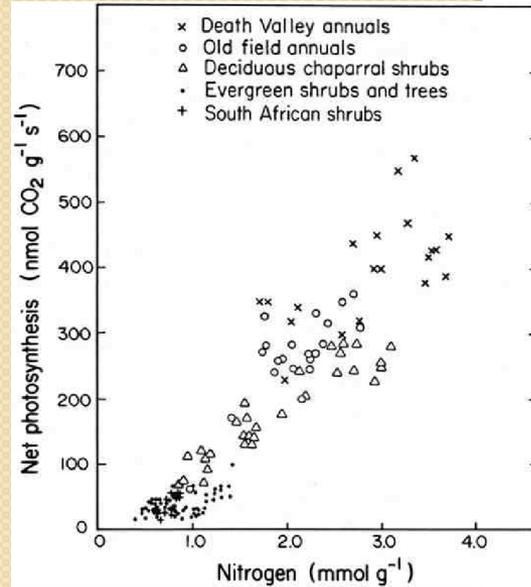
Ex. 2: traits et activité des feuilles



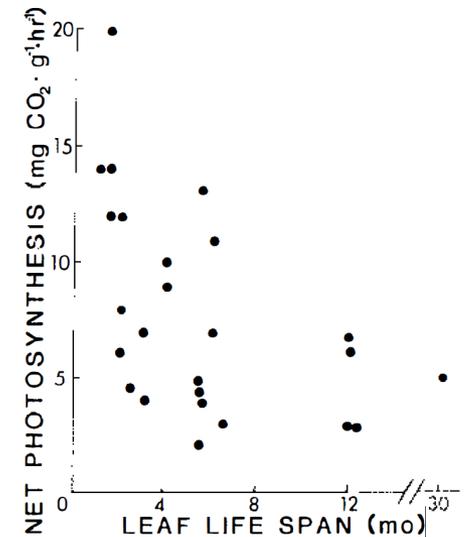
La feuille et la photosynthèse: porte d'entrée dans le monde vivant



QUELS SONT LES CONTRÔLES SUR LA VITESSE DE PHOTOSYNTHÈSE ?

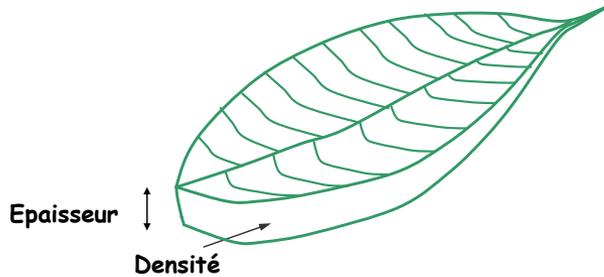


Field & Mooney (1986)
On the Economy of Plant Form



Chabot & Hicks (1982)
Ann Rev Ecol Syst 13: 229

I – La structure morpho-anatomique des feuilles



**SSF (SLA) = Surface d'une
feuille/Masse sèche**

[MSF (LMA) = Masse sèche/surface]

30 espèces d'arbres de
forêt tempérée (Japon)

23 espèces d'arbres de
la forêt amazonienne

29 espèces herbacées
d'Europe centrale

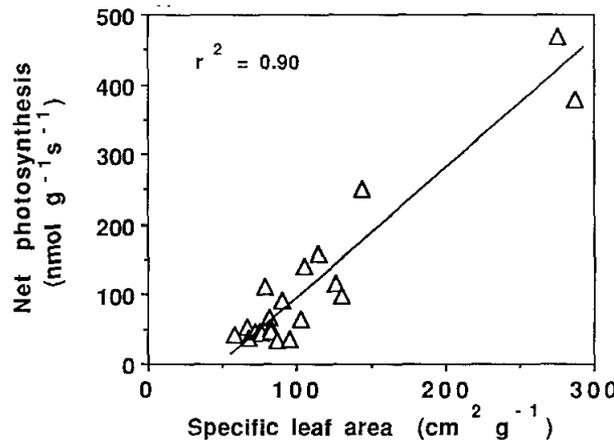
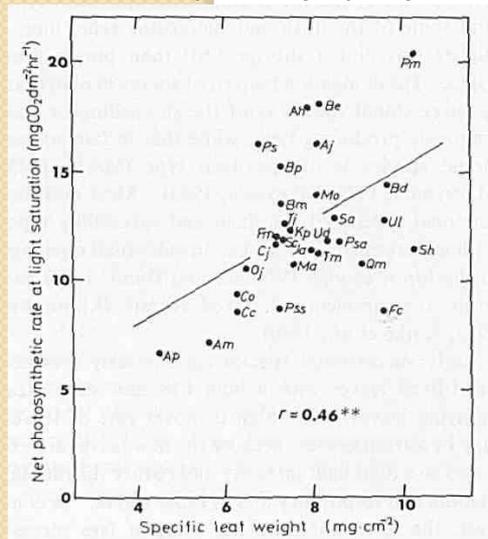


Table 4. Leaf life span and leaf properties for a subsample of the species listed in Table 1 and 2

Plant species	Longevity (d)	LWA (g m ⁻²)	Nitrogen content (% d.m.)	A _{exp} ^a (mmol m ⁻²)	LWA/A _{exp} ^b
Low altitude:					
<i>Am.</i>	42	51.8	3.88	193	17.9
<i>Ra.</i>	55	65.4	2.93	–	16.1
<i>Pa.</i>	57	52.3	3.09	116	20.2
<i>Tr.</i>	61	33.9	4.50	109	16.3
<i>Cc.</i>	61	37.9	3.50	95	19.0
<i>Cl.</i>	66	51.6	3.46	156	22.8
<i>To.</i>	66	34.4	3.60	81	22.7
<i>Cp.</i>	78	32.4	–	–	1.51
<i>Gr.</i>	87	52.6	2.96	91	12.5
<i>Pc.</i>	94	42.3	1.82	54	10.7
High altitude:					
<i>Pv.</i>	41	69.0	4.51	193	20.2
<i>Od.</i>	48	25.4	5.00	129	16.4
<i>Dc.</i>	54	48.8	4.04	141	21.1
<i>La-p.</i>	58	68.0	2.98	145	14.7
<i>Pa.</i>	61	64.5	2.27	136	20.2
<i>Gm.</i>	64	69.9	2.81	175	14.9
<i>Ta.</i>	66	36.5	2.94	77	10.6
<i>Ae.</i>	70	45.0	4.62	168	16.4
<i>Pc.</i>	73	60.6	3.28	168	16.4
<i>Pg.</i>	76	74.1	2.47	139	12.3
<i>Lm.</i>	77	79.4	3.62	200	23.9
<i>Eu.</i>	78	47.6	3.63	138	22.5
<i>La.</i>	80	62.5	2.36	111	11.0
<i>Gr.</i>	91	70.4	2.89	152	13.7
<i>Rg.</i>	93	64.5	3.37	173	19.4

^a A_{exp} μmol m⁻² s⁻¹; ^b LWA/A_{exp} g s μmol⁻¹

* P ≤ 0.05, ** P ≤ 0.01, n.s. not significant

Koike (1988)
Pl Sp Biol 3: 77

Reich *et al.* (1991)
Oecologia 86: 16

Diemer *et al.* (1992)
Oecologia 89: 10

II – La durée de vie des feuilles

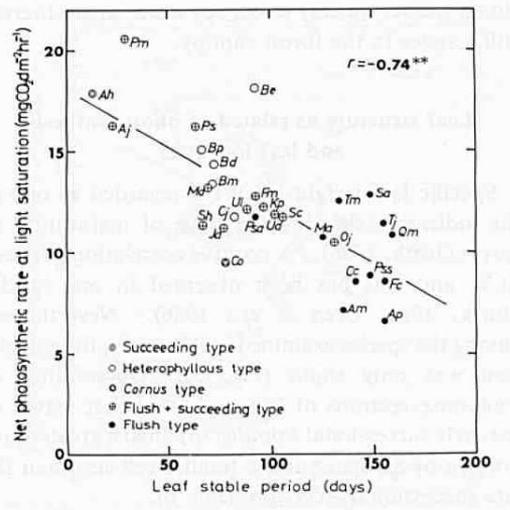


30 espèces d'arbres de forêt tempérée (Japon)

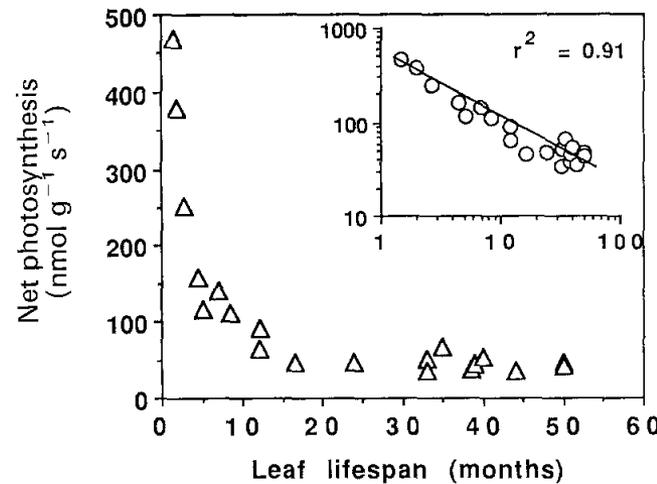
23 espèces d'arbres de la forêt amazonienne



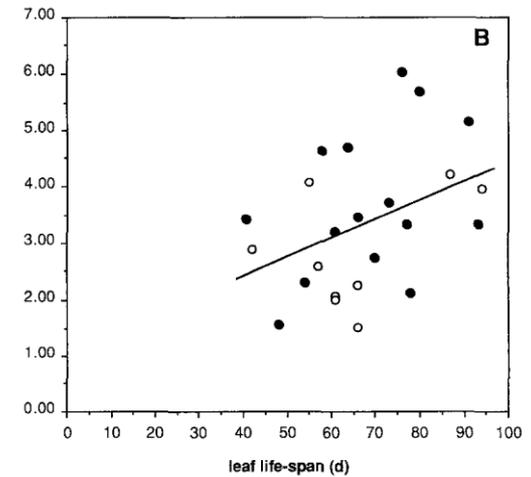
29 espèces herbacées d'Europe centrale



Koike (1988)
Pl Sp Biol 3: 77

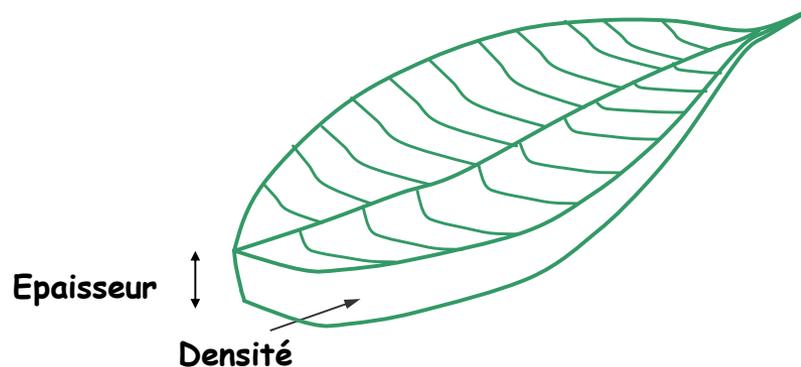


Reich *et al.* (1991)
Oecologia 86: 16



Diemer *et al.* (1992)
Oecologia 89: 10

Homogénéisation des variables pour la synthèse => travail sémantique (I)



$$\begin{aligned} \text{SLA} &= \text{surface/masse} \\ &= 1/\text{LMA (SLW, LWA)} \end{aligned}$$

- Koike (1988) and Diemer *et al.* (1992) presented data on leaf mass per area, whereas Reich *et al.* (1991) provided data on the inverse, SLA. All data were converted to SLA to enable comparison with data from Reich *et al.* (1992), but otherwise, use of this measure, rather than its inverse, is arbitrary. Neither Koike (1988) nor Diemer *et al.* (1992) expressed A_{max} on a mass basis, but they did provide A_{max} values on an area basis, which I combined with SLA data to obtain estimated mass-based A_{max} rates.

Homogénéisation des variables pour la synthèse => travail sémantique (2)



#107929703

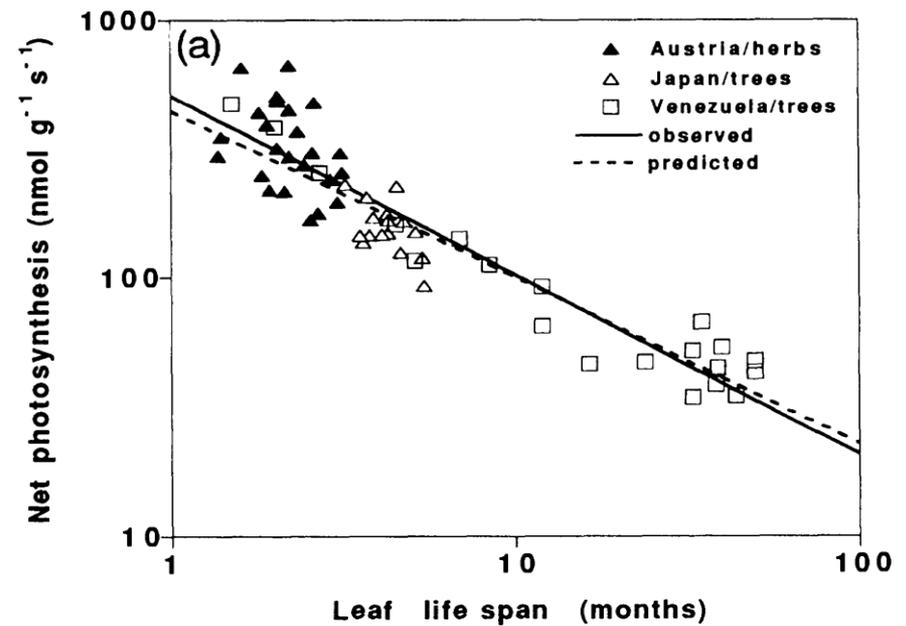
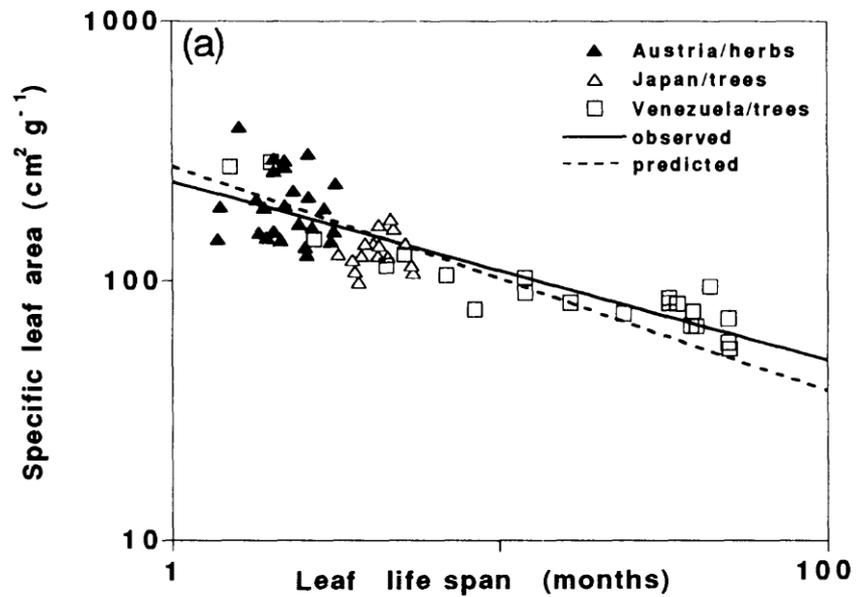


Leaf life span
≠
Leaf stable period
≠
Functional life span

- Diemer *et al.* (1992) and Reich *et al.* (1991) provided data on mean and median leaf life spans, respectively, but in the latter study, mean and median leaf life spans were similar. Diemer *et al.* (1992) used a measure of the 'functional leaf life span' that did not include the time required for the final decay of leaves during senescence prior to shedding. Koike (1988) provided both a measure of the 'leaf stable period' and a measure of the mean leaf life span. The 'leaf stable period' was defined to include only the period when A_{max} was stable and maximal. Thus, the 'leaf stable period' is less than the total leaf life span, and also often less than the 'functional life span' used by Diemer *et al.* (1992). Therefore, mean life-span data (available for 16 species, Koike 1988) were used.

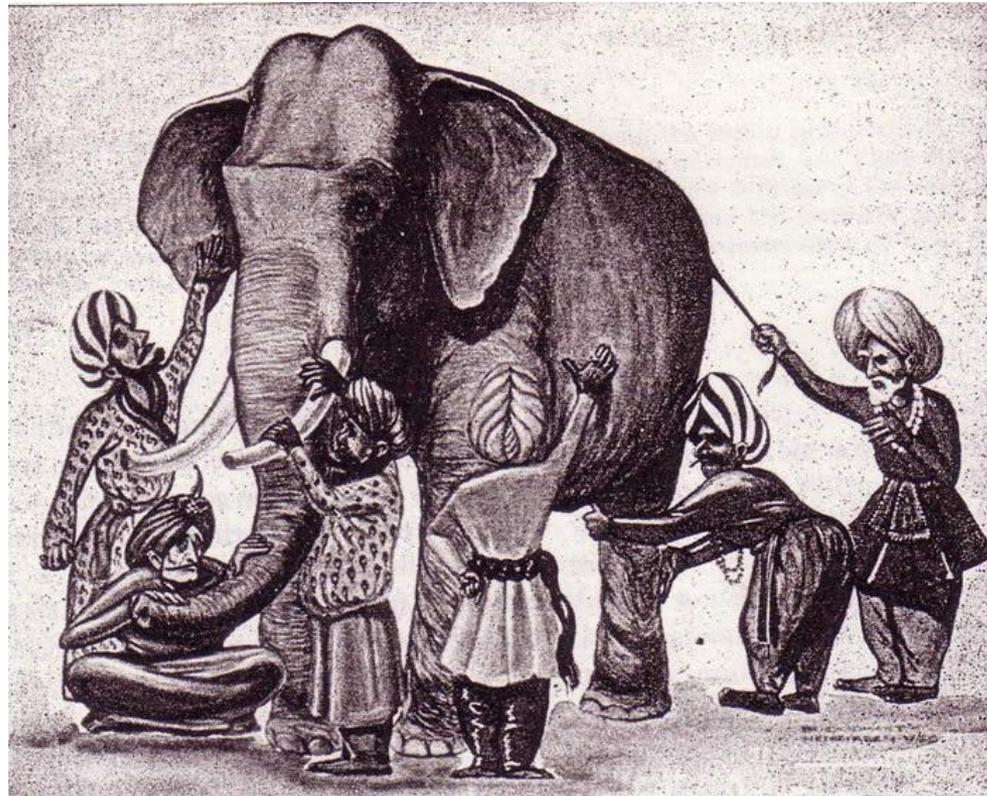
Reich *et al.* (1993)
Functional Ecology 7: 721

Test sur les jeux de données combinés



Reconciling apparent discrepancies among studies relating life span, structure and function of leaves in contrasting plant life forms and climates: ‘the blind men and the elephant retold’

P. B. REICH



La parabole des aveugles (*i.e.* les scientifiques) et de l'éléphant (*i.e.* les lois générales)

Le réseau « Global Plant Trait Network »: 40 chercheurs, 15 pays

Fonction d'acquisition du carbone:

- vitesses de photosynthèse max. et de respiration
- surface spécifique et durée de vie
- composition chimique: N et P

Base de données mondiale:

- 2548 espèces (~1% du nombre total d'espèces de végétaux terrestres);
- 219 familles; aucune restriction sur le type d'espèce (fougères, herbacées, arbustes, arbres)

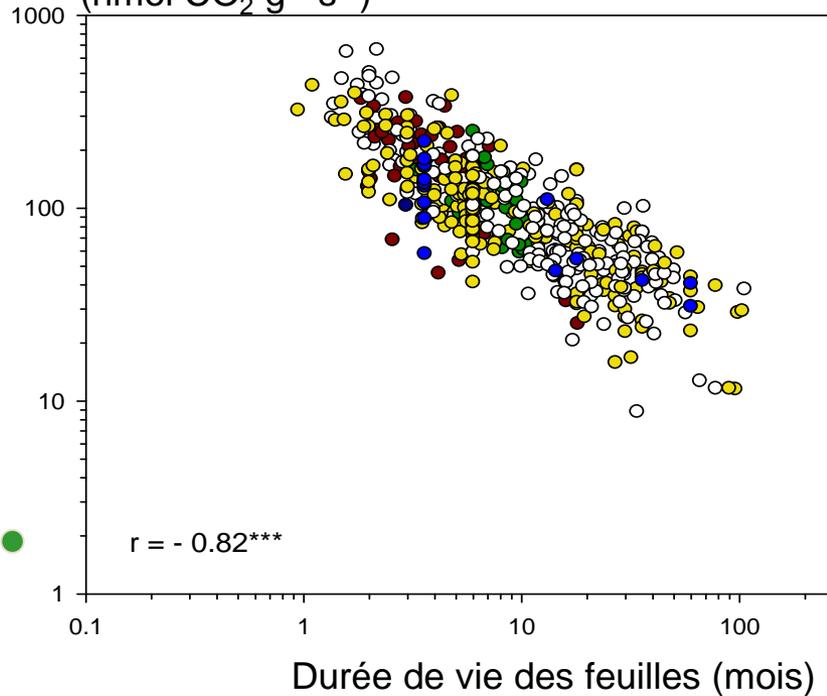


175 sites

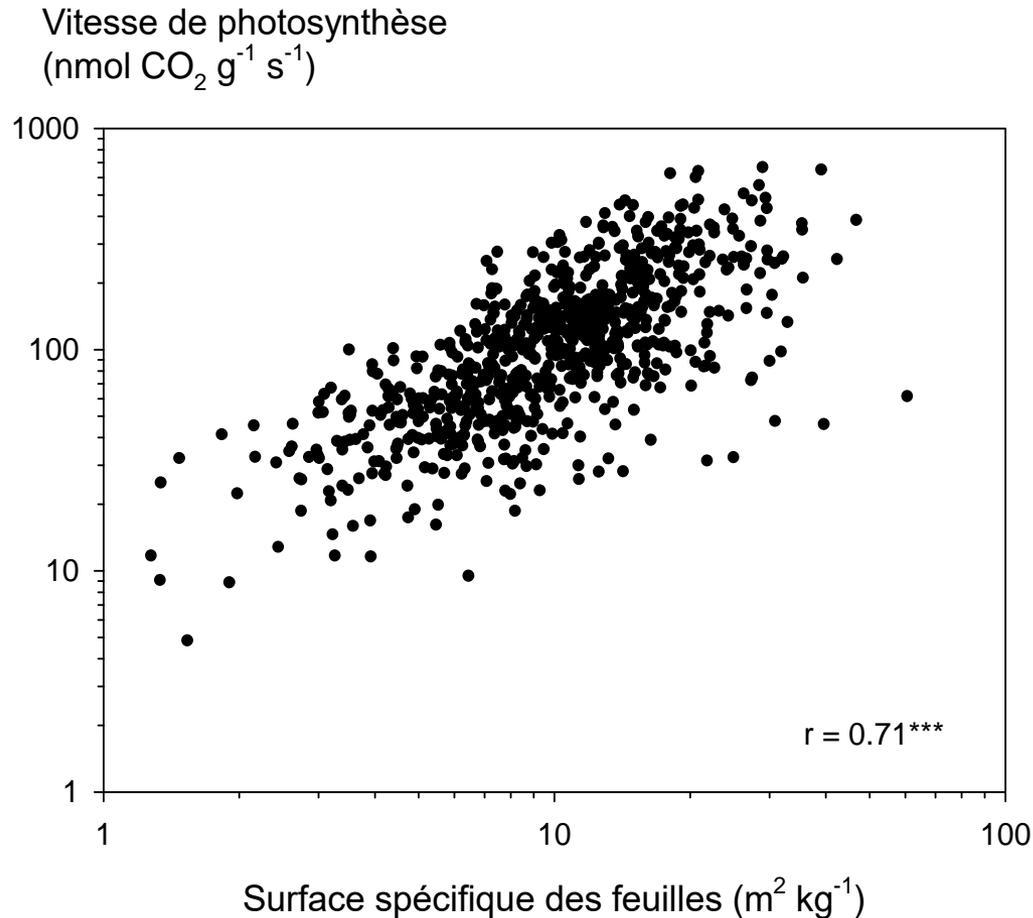
Un schéma universel de fonctionnement des végétaux (I)

- Autres biomes

Vitesse de photosynthèse
($\text{nmol CO}_2 \text{ g}^{-1} \text{ s}^{-1}$)



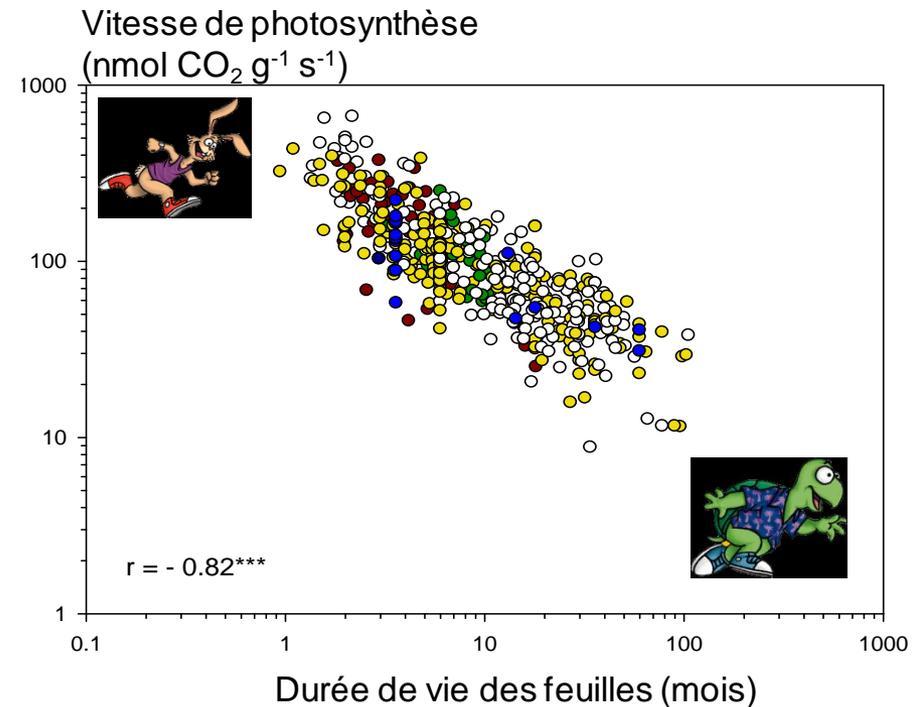
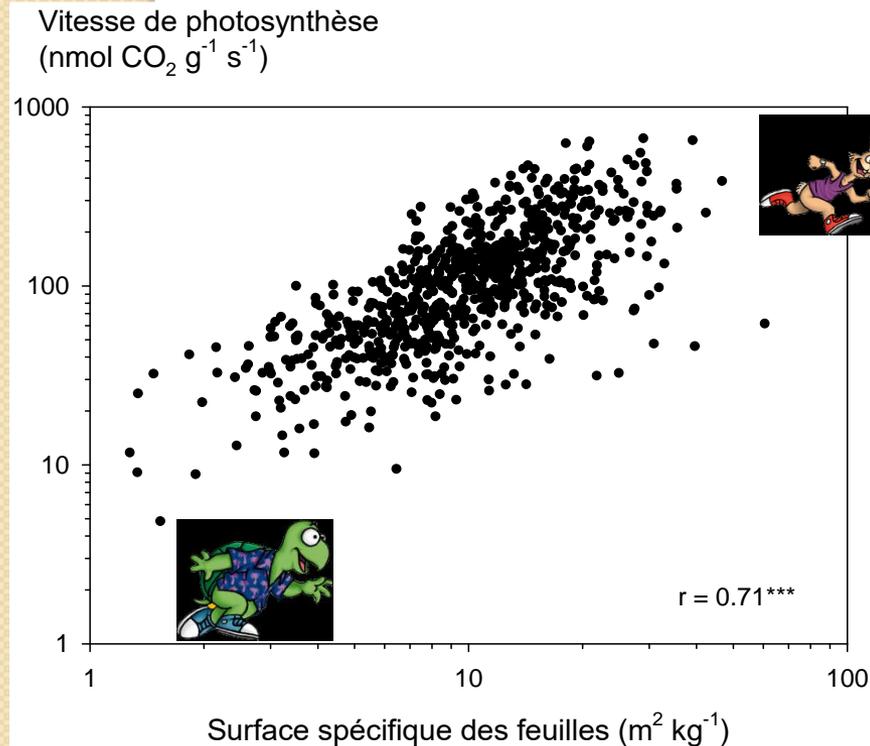
Un schéma universel de fonctionnement des végétaux (2)



The worldwide leaf economics spectrum

Ian J. Wright¹, Peter B. Reich², Mark Westoby¹, David D. Ackerly³, Zdravko Baruch⁴, Frans Bongers⁵, Jeannine Cavender-Bares⁶, Terry Chapin⁷, Johannes H. C. Cornelissen⁸, Matthias Diemer⁹, Jaume Flexas¹⁰, Eric Garnier¹¹, Philip K. Groom¹², Javier Gulias¹⁰, Kouki Hikosaka¹³, Byron B. Lamont¹², Tali Lee¹⁴, William Lee¹⁵, Christopher Lusk¹⁶, Jeremy J. Midgley¹⁷, Marie-Laure Navas¹¹, Ülo Niinemets¹⁸, Jacek Oleksyn^{2,19}, Noriyuki Osada²⁰, Hendrik Poorter²¹, Pieter Poort²², Lynda Prior²³, Vladimir I. Pyankov²⁴, Catherine Roumet¹¹, Sean C. Thomas²⁵, Mark G. Tjoelker²⁶, Erik J. Veneklaas²² & Rafael Villar²⁷

NATURE | VOL 428 | 22 APRIL 2004 | www.nature.com/nature



Les leçons à tirer de cette course d'obstacles

- Les données doivent pouvoir être découvertes
- Leur réutilisation (par la personne qui les a collectées ou par des tiers) nécessite une mise en forme soignée et explicite
- C'est une histoire qui finit (plutôt) bien: mise en évidence d'une loi fondamentale du fonctionnement des plantes
- Mais... que d'efforts !

Vers une standardisation...

CSIRO PUBLISHING

www.publish.csiro.au/journals/ajb

Australian Journal of Botany, 2003, 51, 335–380

A handbook of protocols for standardised and easy measurement of plant functional traits worldwide

*J. H. C. Cornelissen^{A,J}, S. Lavorel^B, E. Garnier^B, S. Díaz^C, N. Buchmann^D, D. E. Gurvich^C,
P. B. Reich^E, H. ter Steege^F, H. D. Morgan^G, M. G. A. van der Heijden^A,
J. G. Pausas^H and H. Poorter^I*

 Global Change Biology

Global Change Biology (2011) 17, 2905–2935, doi: 10.1111/j.1365-2486.2011.02451.x

TRY – a global database of plant traits

J. KATTGE*, S. DÍAZ†, S. LAVOREL‡, I. C. PRENTICES§, P. LEADLEY¶, G. BÖNISCH*,
E. GARNIER||, M. WESTOBY§, P. B. REICH**, ††, I. J. WRIGHT§, J. H. C. CORNELISSEN ‡‡

Journal of Ecology



Journal of Ecology 2017, 105, 298–309

doi: 10.1111/1365-2745.12698

Towards a thesaurus of plant characteristics: an ecological contribution

Eric Garnier^{*1,2}, Ulrike Stahl^{3,4,5}, Marie-Angélique Laporte^{1,4,6}, Jens Kattge^{3,4},
Isabelle Mougenot⁷, Ingolf Kühn^{4,5,8}, Baptiste Laporte², Bernard Amiaud^{9,10},
Farshid S. Ahrestani^{11,12}, Gerhard Bönisch³, Daniel E. Bunker¹³, J. Hans C. Cornelissen¹⁴,
Sandra Díaz¹⁵, Brian J. Enquist¹⁶, Sophie Gachet¹⁷, Pedro Jaureguiberry¹⁵,
Michael Kleyer¹⁸, Sandra Lavorel¹⁹, Lutz Maicher^{20,21}, Natalia Pérez-Harguindeguy¹⁵,
Hendrik Poorter²², Mark Schildhauer²³, Bill Shipley²⁴, Cyrille Violle¹, Evan Weiher²⁵,
Christian Wirth^{4,26}, Ian J. Wright²⁷ and Stefan Klotz⁵

Plus de détails demain!

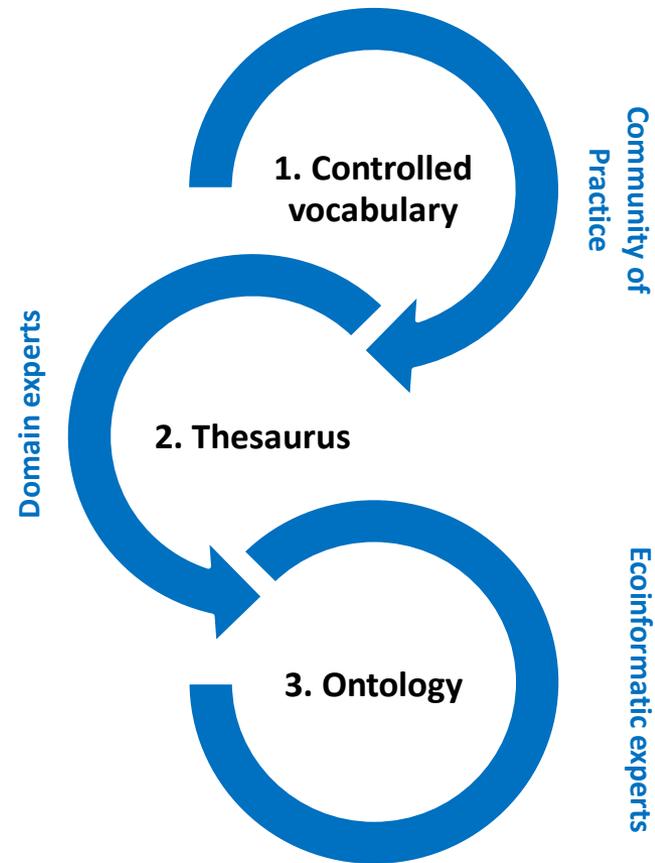
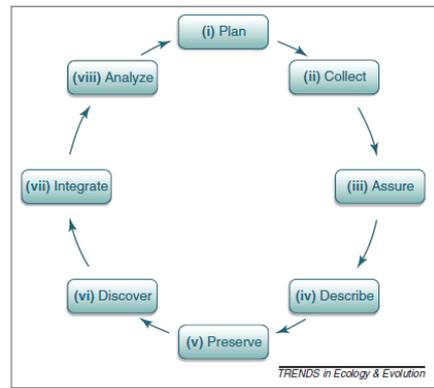
MERCI POUR VOTRE ATTENTION

DES QUESTIONS?



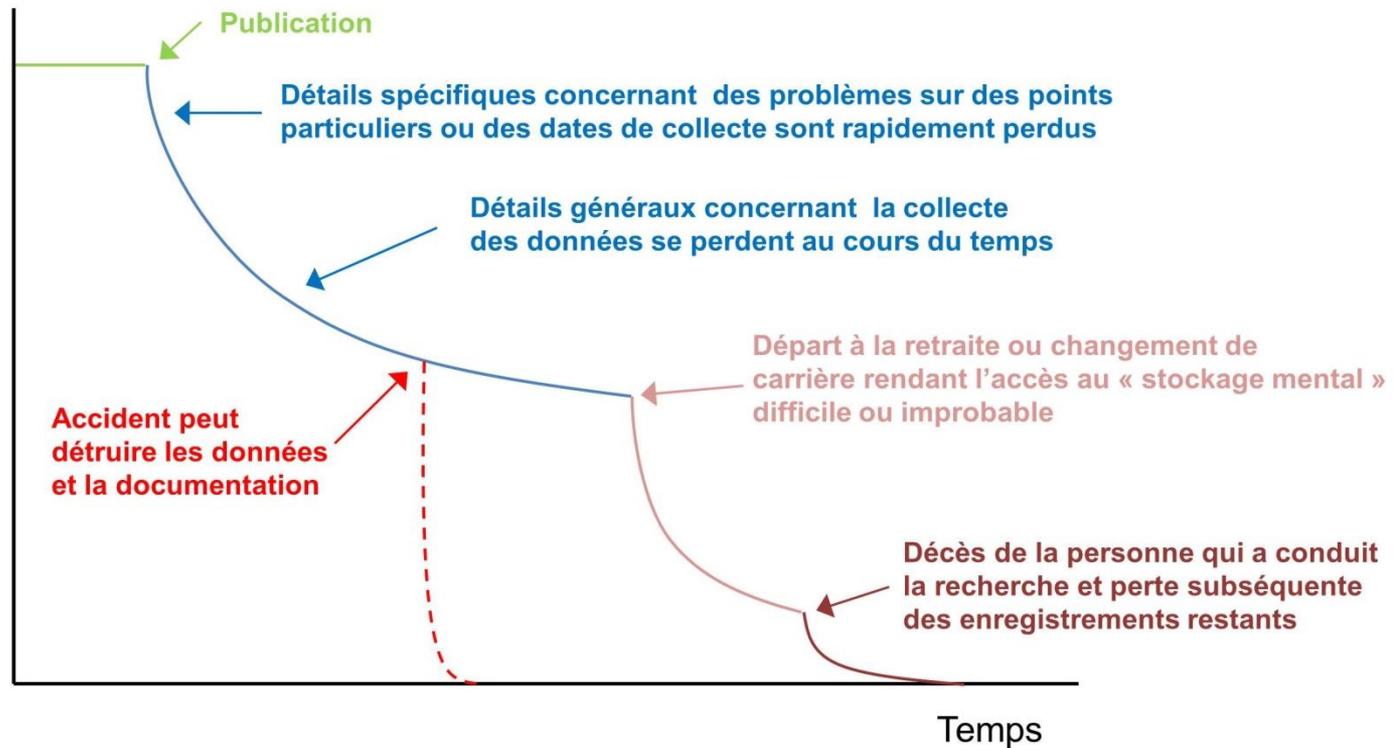


Cycle de vie des données: (VII) intégrer: sémantique

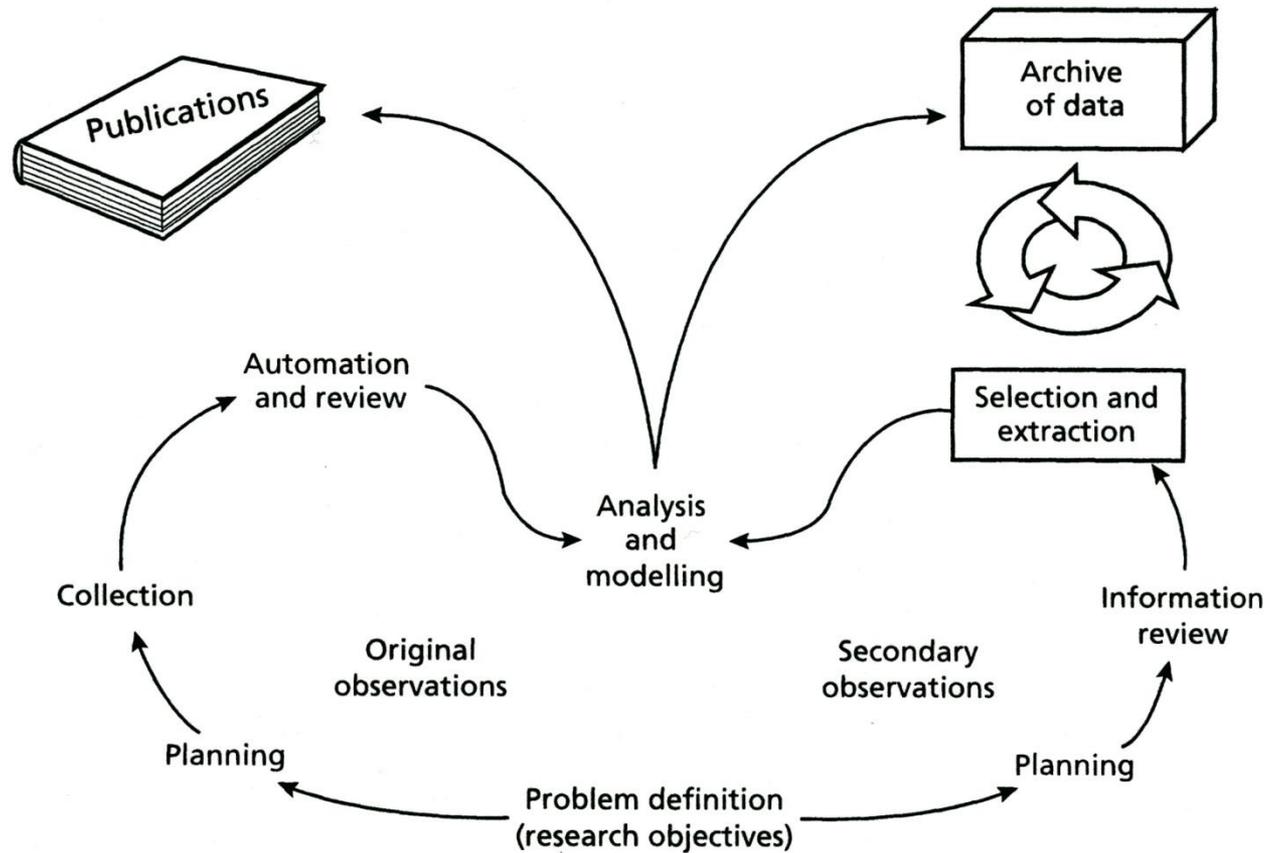


Données et information: une histoire d'entropie

Contenu en information des
données et des métadonnées



Les données et la recherche



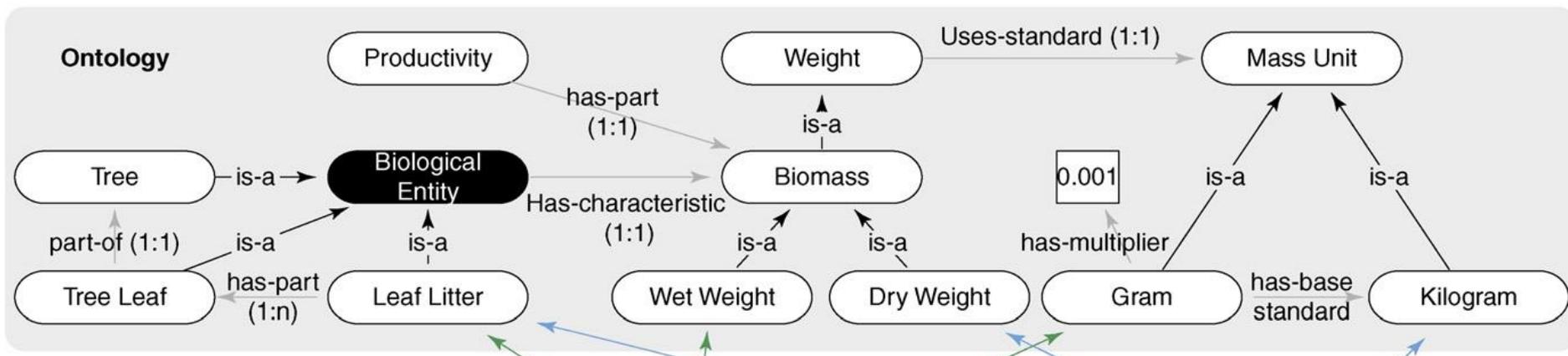
Intégrer: les ontologies

- En sciences de l'information, une ontologie est un mode formel d'organisation de la connaissance dans lequel:
 - chaque élément (ou concept) est défini précisément
 - chaque relation potentielle entre ces éléments est paramétrée ou contrainte (e.g. is_a, part_of, has_member, has_characteristic...)

Schuurman & Leszczynski (2008) *BBI* 2: 187

- Représentation explicite d'un domaine permettant d'effectuer des tâches raisonnées de façon automatique

Intégrer: exemple d'ontologie en écologie



Pour éviter ça...



Adéquation entre utilisation et effort de structuration

Level	Planned Use			
III	Publishable and auditable	Inadequate	Minimal	Good practice
II	Searchable and third party reuse	Minimal	Good practice	Excessive
I	Exchange with expert colleague	Good practice	Excessive	Excessive
		LOW	MEDIUM	HIGH

Michener *et al.* (1997)
Ecol Appl 7 : 330

Amount of structure
{Formalization, level of effort}

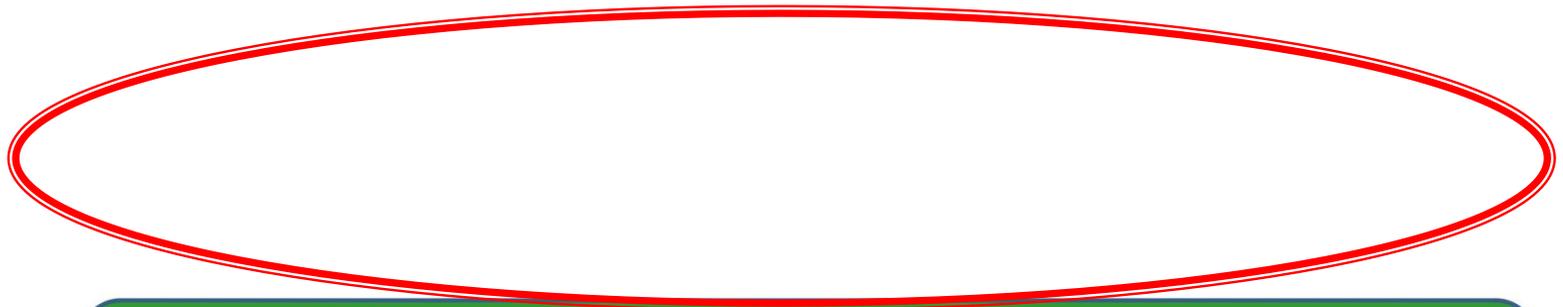
En guise de conclusion (2)

- Les données collectées peuvent être utilisées pour des questions qui ne sont pas encore formulées
- Pensez plus loin que votre propre étude pour que vos données puissent être réutilisées
- En route pour le « big data » en écologie... mais pas seulement!

L'émergence de l'écoinformatique:

- champ de recherche et de développement situé aux interfaces de l'écologie, de l'informatique et des sciences de l'information...
(Jones *et al.* 2006, *Ann Rev Ecol Syst*)
- ... qui a pour objectif de permettre aux scientifiques de générer de nouvelles connaissances grâce à des outils et approches innovants pour découvrir, gérer, intégrer, analyser, visualiser et sauvegarder les données et les informations biologiques, environnementales et socioéconomiques pertinentes
(Michener & Jones 2012, *TREE*).

Les objectifs de l'écoinformatique



(c) Outils pour les scientifiques et les gestionnaires de données



Publication des données

- annotation sémantique
- automatisée (si possible)

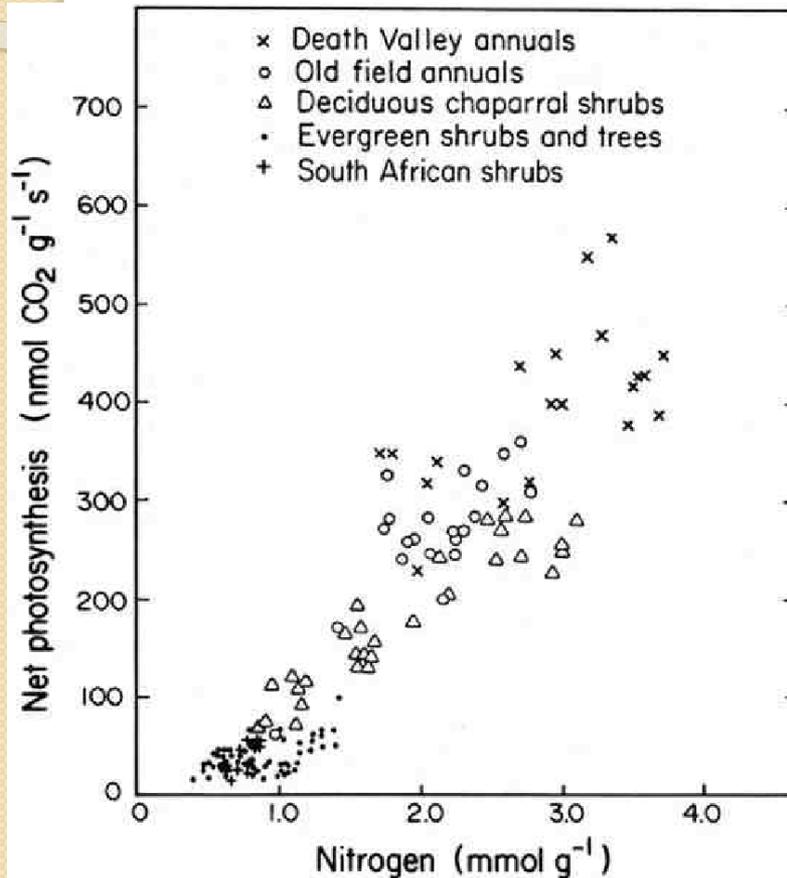
Découverte des données

- fondée sur les ontologies
- résultats résumés/triés

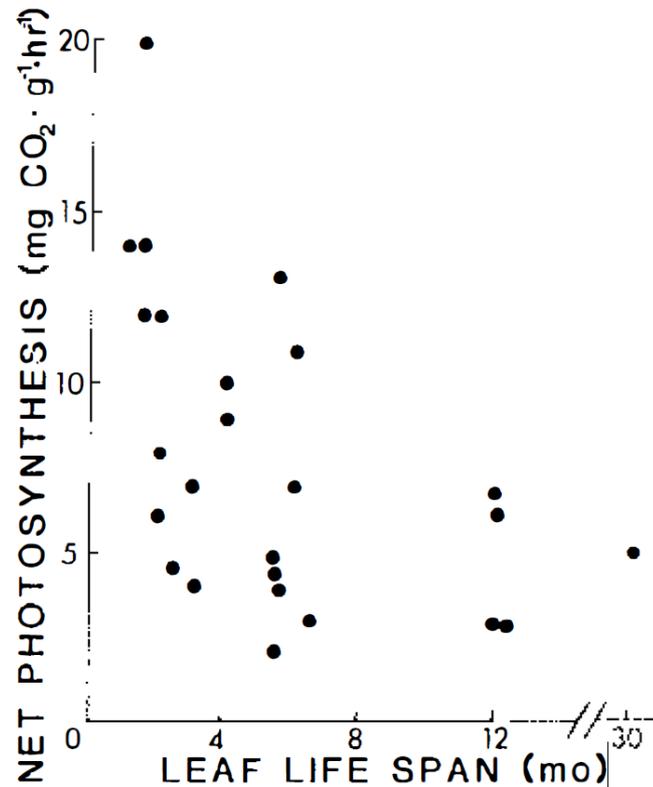
« Workflows » scientifiques

- accéder et intégrer les données
- créer, conduire des analyses

Quels sont les contrôles sur la vitesse de photosynthèse?



Field & Mooney (1986)
On the Economy of Plant Form



Chabot & Hicks (1982)
Ann Rev Ecol Syst 13: 229