

Authors: Romain David¹, Alison Specht², Margaret O'Brien³, Lesley Wyborn⁴, Christina Drummond⁵, Rorie Edmunds⁶, Claudia Filippone¹, Jeaneth Machicao⁷, Nobuko Miyairi⁸, Graham Parton⁹, Debora Pignatari Drucker¹⁰, Shelley Stall¹¹, Niklas Zimmer¹²

Contact: R. David romain.david@erinha.eu

[@Romain_DAVID_13](https://twitter.com/Romain_DAVID_13)

Websites: <https://parsecproject.org/> <https://www.erinha.eu/>

twitter: @PARSEC_News @ERINHA_RI

MULTILINGUAL CHALLENGES

Episode 0
PROFESSIONALIZING
DATA STEWARDSHIP AT THE
GLOBAL SCALE

Profound changes in our world are exacerbating data availability challenges at the global level, in particular between scientists and other knowledge workers from regions separated by various features including historical, financial, cultural, political aspects, aside from time and space ...

Very few, if any, of our present problems such as biodiversity decline, climate change, and viral pandemics stop at national, disciplinary and linguistic boundaries, yet our most vital responses to the shared problems, the information generated to analyze and derive solutions, **is still siloed in different languages and locations throughout the world.** It is clear that in order for us to effectively respond, we need to collaborate globally and communicate information more effectively. Globalization of research requires interoperability of our observations and experimentation systems.

My name is **C-3PO!**
I am fluent in over six million forms of communication... **And you?**



Good practices

For translators:

- Scientific skills with at least good command of target languages
- Validators to check work translated
- A clear versioning system
- A translating strategy with prioritization levels

For Databases:

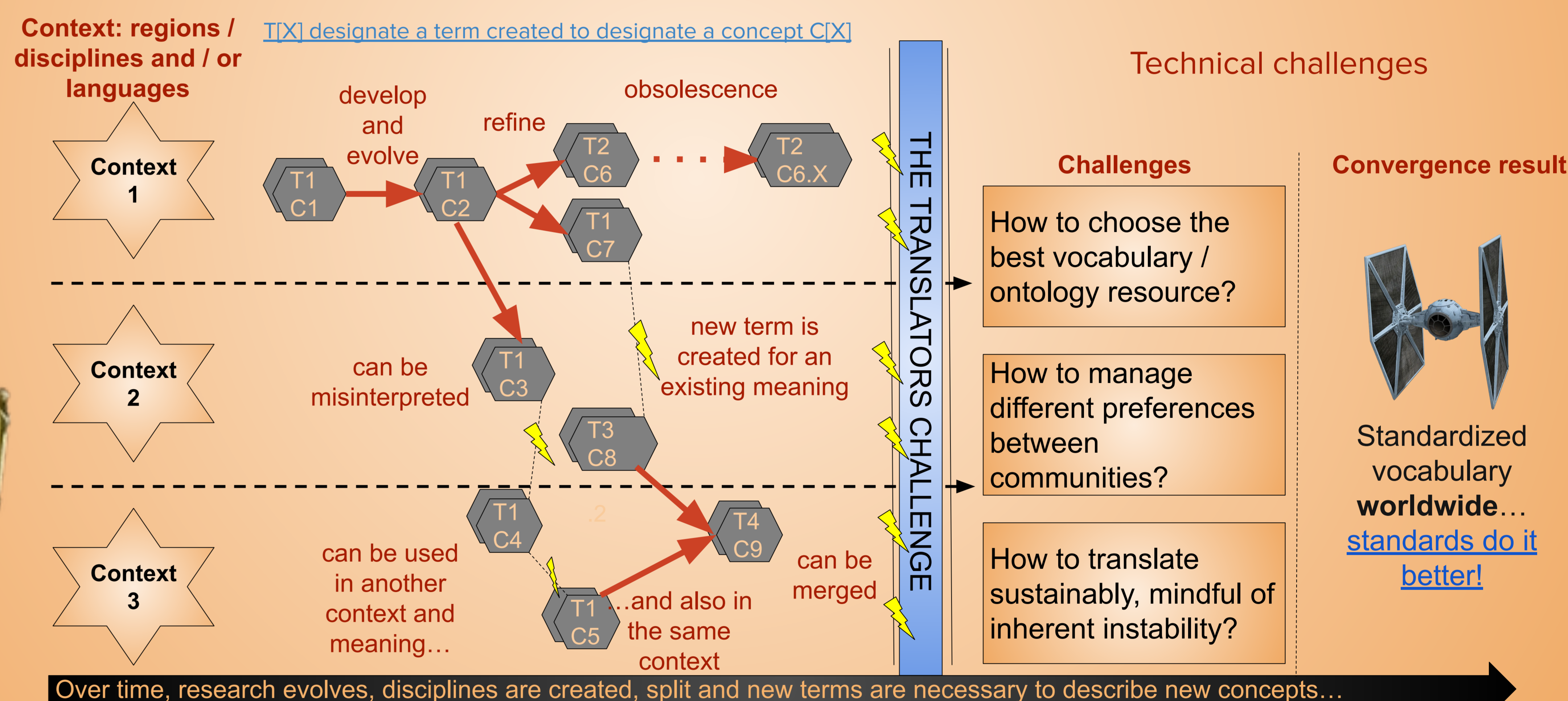
- Data dictionary with
 - disambiguated definitions both linked to and linked by concepts in shared ontologies
 - the date of the definition
 - the context of use (disciplinary - protocol - necessary skills...)
- A description of the community approval process (that can be adopted for new terms by other communities)

Science must be universally shared, in an interdisciplinary way, respecting that all disciplines are creating new concepts. The fact that science is always creating them is the core challenge for translators

The "TERM WARS" at global Scale: a time and space story

Discrepancies between regions and groups (culture, content, workflows, language, semantics, translation, funds...) are numerous whether for individual data users or more universally. We must anticipate issues such as **which is the preferred language, polysemy** (1 term, multiple meanings), **confusion** (multiple terms for 1 meaning or 'false friends' between 2 languages), plus **existing and evolving nuances** (not an exact match between languages and during time). Furthermore, terms are often **adopted from another language with different contexts and disciplinary realms** (that might decrease interoperability) and impedes **translation of all versions at the same time.**

Translation occurs at the concept level, *not as a simple one to one translation of (consecutive) words.*



Over time, research evolves, disciplines are created, split and new terms are necessary to describe new concepts...

Translation challenges are also organisational and sustainability challenges:

Care must be taken to ensure that datasets resulting from projects that practice co-creation and co-evolution of knowledge are translated into indigenous languages such that they can be used by the affected communities. Taking these challenges into account, we have to consider **human effort and the level of translation**, e.g. a low or minimum yet sustainable level, that is *legally allowable*. How can such minimal objectives be linked with FAIR principle compliance (especially FAIR and community approved vocabularies)? In several communities, translation is voluntary. One of the sustainability challenges is the need for ongoing engagement: how to keep interested groups involved. We need expert translators, as described, to **maintain the quality level critical to achieve effective harmonization** among languages.



Acknowledgements:
PARSEC is funded by the Belmont Forum through the National Science Foundation (NSF), The São Paulo Research Foundation (FAPESP), the French National Research Agency (ANR), and the Japan Science and Technology Agency (JST). This work is partially funded by the EOSC-Life European program (grant agreement No. 824087). We acknowledge wikipedia for R2D2, C-3PO and Tie fighter pictures.

Author affiliations : ¹ERINHA (European Research Infrastructure on Highly Pathogenic Agents) AISBL, FR, ORCID: 0000-0003-4073-7456; ²The University of Queensland, AU, ORCID: 0000-0002-2623-0854; ³University of California Santa Barbara, USA, ORCID:0000-0002-1693-8322; ⁴Australian National University, AU, ORCID: 0000-0001-5976-4943; ⁵EducoPIA Institute, ORCID: 0000-0001-5794-0413; ⁶DataCite; ⁷University of São Paulo, ORCID: 0000-0002-1202-0194, ⁸Nobuko Miyairi, ORCID: 0000-0002-3229-5662, ⁹Center for Environmental Data Analysis, ORCID: 0000-0003-4157-0352, ¹⁰Embrapa Digital Agriculture, ORCID: 0000-0003-4177-1322, ¹¹American Geophysical Union, ORCID: 0000-0003-2926-8353, ¹²University of Cape Town ORCID: 0000-0001-8078-0403