# Graded resources:
# from linguistic engineering
# to practical applications

**Núria Gala**

December 7th, 2021 - Louvain-la-Neuve

# Overview

*Graded resources: roots and context*

Back to the roots : empiricism

Grade schools, reading standards

Readability, word-lists and corpus linguistics

Graded resources



*Linguistic engineering and practical applications*

Using graded reading materials

Analyzing language complexity

Modelling student difficulties

Using the resources in the classrooms or in personalized training

# *Graded resources: roots and context*

Back to the roots : empiricism

Grade schools, reading standards

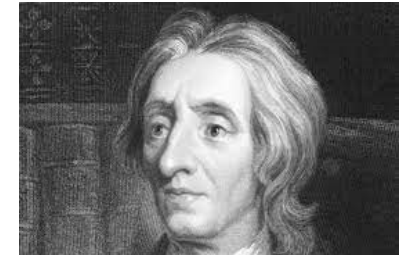Readability, word-lists and corpus linguistics

Graded resources

# Empiricism

Derivation from the ancient Greek word *empeiria* ("experience")

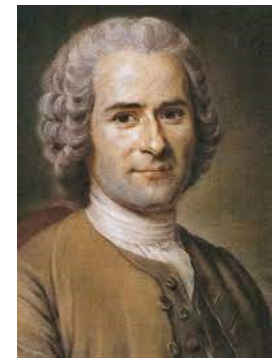John Locke (1632 – 1704)
British empiricist: "truth and knowledge arise out of observation and experience rather than manipulation of accepted or given ideas". Need for children to have concrete experiences to learn.

Jean Jacques Rousseau (1712 – 1778)
His philosophy of education: learning through experiencing, "child-centered" education.

Influence on **education**: grade schools, reading tests, adapted reading materials. **Putting the learner at the forefront.**
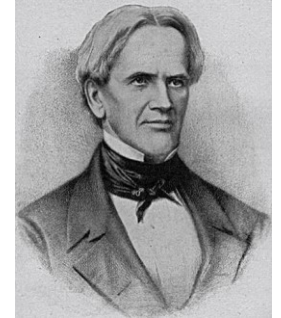
# From one-room schools... to grade schools



Early 19th « Children under the age of five were often mixed in with adults in their twenties. Additionally, classrooms were frequently overcrowded, housing as many as eighty students at a time. Because of the overcrowding, already scarce textbooks and learning materials had to be spread even more thinly amongst students. As a result, class time amounted to a tedious recitation of facts and instructors struggled to devote individual attention to students. »

Ted Brackemyre. *The Rise of Public Education in Early America*. 2021 U.S. History Scene.
https://ushistoryscene.com/article/rise-of-public-education/

Horace Mann (1796 – 1859), promoter of public education.

1847 first graded school (Boston) with books prepared for each grade.



Students learn best with materials written for their current reading level.

Reading standards were set for each grade.

| Multi-aged one-room schools. Scarce non adapted textbooks. | → | Students grouped by grades. Standards adapted to each grade. |
|---|---|---|

# Empiricism

Derivation from the ancient Greek word *empeiria* ("experience")

L. Bloomfield and Z. Harris, US distributionalism (1940s-1960s), based on behaviorist psychological theories and on direct observation of environments: *"you shall know a word by the company it keeps"* (Firth, 1957).

J. Sinclair and G. Leech (1970s-2000s), Corpus linguistics: study of language through its samples, e.g., corpus-driven lexicons for foreign learners of English.

Influence on **language teaching**: readability formulae, word-lists, corpus.
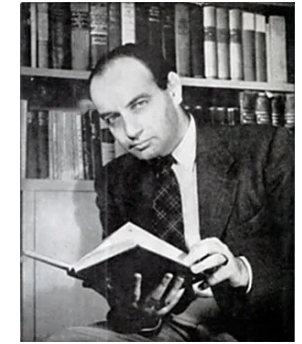**Putting word distributions at the forefront.**

# Predicting text readability
# Vocabulary learning through word lists

First readability formulae (20th)

- B. A. Lively and S. L. Pressey, predicting readability based on word-frequencies (*A method for measuring the vocabulary burden of text-book*, 1923)
- **R. Flesch** (*Marks of Readable Style : A study in Adult Education,* 1943 et 1948)
  Reading Ease score: length (syllables/word, words/sentence)

Computational readability (21st)

NLP and machine learning Collins-Thompson & Callan (2005), François (2009)

Neural approaches, deep learning Deutsch, Jasbi & Shieber (2020), Martinc, Pollak & Robnik-Šikonja (2021)

*Teachers' Book of Words* (Thorndike, 1921)
*Basic English* (Ogden, 1930)

**a to acacia**

| | G | T | L | J | S | | G | T | L | J | S |
|---|---|---|---|---|---|---|---|---|---|---|---|
| a | AA | M | M | M | M | aborigines | 1 | 7 | 8 | 5 | 12 |
| Aaron | 2 | 28 | 6 | 5 | 14 | abortive | 1 | 11 | 1 | 3 | 15 |
| aback | 2 | 10 | 15 | 11 | 12 | abound | 12 | 90 | 32 | 39 | 59 |
| abandon | 38 | 119 | 150 | 130 | 285 | about | AA | M | M | M | M |
| abandoned (adj.) | 3 | 11 | 14 | 12 | 27 | above | AA | M | 941 | M* | ? |
| abandonment | 3 | 10 | 16 | 3 | 39 | Abraham | 11 | 115 | 47 | 26 | 22 |
| abase | 1 | 14 | 2 | 0 | 5 | Abram | 1 | 7 | 0 | 0 | 14 |
| abash | 3 | 16 | 14 | 24 | 13 | abreast | 4 | 16 | 17 | 23 | 20 |
| abate | 7 | 57 | 20 | 20 | 33 | abridge | 2 | 18 | 0 | 6 | 13 |
| abatement | 1 | 10 | 5 | 2 | 4 | abridgment | 1 | 11 | 1 | 0 | 9 |
| abbé | 3 | 7 | 18 | 0 | 44 | abroad | 48 | 200 | 198 | 200* | 268 |
| abbess | 1 | 14 | 3 | 9 | 1 | abrogate | 1 | 10 | 0 | 2 | 9 |
| abbey | 11 | 57 | 19 | 51 | 83 | abrupt | 6 | 27* | 43 | 20 | 26 |

Thorndike & Lorge (1921)

# Graded resources

Influence on **education**: grade schools, reading tests, adapted reading materials. **Putting the learner at the forefront.**

Influence on **language teaching**: readability formulae, word-lists, corpus. **Putting word distributions at the forefront.**

**Graded resources**: structured series of linguistic data scaled according to the ease (or difficulty) of learning, reading and comprehending.

Remarks:
- Lexicons *vs* corpora
- Scales often correspond to stablished learning grades, i.e., CEFR
- Teacher judgments of the difficulty *vs l*earner abilities

# CEFR (2001)

| | | |
|---|---|---|
| **PROFICIENT USER** | C2 | Can understand with ease virtually everything heard or read. Can summarise information from different spoken and written sources, reconstructing arguments and accounts in a coherent presentation. Can express him/herself spontaneously, very fluently and precisely, differentiating finer shades of meaning even in more complex situations. |
| | C1 | Can understand a wide range of demanding, longer texts, and recognise implicit meaning. Can express him/herself fluently and spontaneously without much obvious searching for expressions. Can use language flexibly and effectively for social, academic and professional purposes. Can produce clear, well-structured, detailed text on complex subjects, showing controlled use of organisational patterns, connectors and cohesive devices. |
| **INDEPENDENT USER** | B2 | Can understand the main ideas of complex text on both concrete and abstract topics, including technical discussions in his/her field of specialisation. Can interact with a degree of fluency and spontaneity that makes regular interaction with native speakers quite possible without strain for either party. Can produce clear, detailed text on a wide range of subjects and explain a viewpoint on a topical issue giving the advantages and disadvantages of various options. |
| | B1 | Can understand the main points of clear standard input on familiar matters regularly encountered in work, school, leisure, etc. Can deal with most situations likely to arise whilst travelling in an area where the language is spoken.  Can produce simple connected text on topics which are familiar or of personal interest. Can describe experiences and events, dreams, hopes & ambitions and briefly give reasons and explanations for opinions and plans. |
| **BASIC USER** | A2 | Can understand sentences and frequently used expressions related to areas of most immediate relevance (e.g. very basic personal and family information, shopping, local geography, employment). Can communicate in simple and routine tasks requiring a simple and direct exchange of information on familiar and routine matters.  Can describe in simple terms aspects of his/her background, immediate environment and matters in areas of immediate need. |
| | A1 | Can understand and use familiar everyday expressions and very basic phrases aimed at the satisfaction of needs of a concrete type. Can introduce him/herself and others and can ask and answer questions about personal details such as where he/she lives, people he/she knows and things he/she has. Can interact in a simple way provided the other person talks slowly and clearly and is prepared to help. |

« The data in the scaling studies were intuitive teacher judgments rather than samples of performance. » (Fulcher, 2010)

« The lack of systematicity may be indicative of some **incongruities in the way reading materials were graded with CEFR levels**, which may call for a more critical reflection.  » (Tack, 2021)

# Linguistic engineering and practical applications

Using graded reading materials

Analyzing language complexity

Modelling student difficulties

Using the resources in the classrooms or in personalized training

# Using graded reading materials

Lété, Sprenger-Charolles & Colé (2004)

Scarce tools for studying child language development.
Frequency effect : one of the earliest empirical observations in cognitive psychology.

## MANULEX

First grade-level lexical database built from text-books of year 2000.

Frequency distributions of words observed across text-books for French L1.

5 grades in primary schools, grouped into 3 (according to the ease of reading): CP 6, CE1 7, CE2 to CM2 8-10 years old.

| Lemme | NLET | SYNT | CP (6 years old) | CE1 (7 years old) | CE2-CM2 (8-10 years old) | CP-CM2 |
|---|---|---|---|---|---|---|
| à | 1 | PRE | 14.660,67 | 14.815,37 | 16.868,45 | 15.846,63 |
| à cloche-pied | 13 | ADV | 0,30 | 5,03 | 0,03 | 1,25 |
| à contrecœur | 12 | ADV | | | 1,27 | 0,77 |
| à croupetons | 12 | ADV | | | 0,03 | 0,02 |
| à jeun | 6 | ADV | | | 2,32 | 1,41 |
| à la saint-glinglin | 19 | ADV | | | 1,54 | 0,92 |
| à l'aveuglette | 14 | ADV | | 0,36 | 0,03 | 0,20 |
| à l'improviste | 14 | ADV | | 0,19 | | 0,01 |
| à mi-course | 11 | ADV | | | 0,34 | 0,20 |
| à rebrousse-poil | 16 | ADV | | | 0,09 | 0,05 |
| à tâtons | 8 | ADV | 0,61 | | 5,02 | 3,51 |
| à tire-d'aile | 13 | ADV | 1,05 | 0,25 | 2,39 | 2,46 |
| à tue-tête | 10 | ADV | 8,26 | 6,14 | 5,08 | 7,04 |
| à vau-l'eau | 11 | ADV | | | 0,05 | 0,03 |
| aardvark | 8 | NP | | 1,45 | | 0,05 |
| abaissé | 7 | ADJ | | | 0,36 | 0,21 |
| abaisser | 8 | VER | | 4,16 | 10,24 | 7,83 |
| abajoue | 7 | NC | | | 0,02 | 0,01 |
| abandon | 7 | NC | | | 2,38 | 1,44 |
| abandonné | 9 | ADJ | 4,02 | 17,43 | 15,82 | 15,39 |
| abandonner | 10 | VER | 43,09 | 56,64 | 99,31 | 84,50 |
| abasourdi | 9 | ADJ | | 0,17 | 4,09 | 2,91 |
| abasourdir | 10 | VER | | 0,17 | 0,02 | 0,19 |

http://manulex.org

# FLELex: grading French L2 vocabulary

François, Gala, Watrin & Fairon (2014) ; Tack, François & Fairon (2016)

Word frequencies by difficulty level according the CEFR scale, first resource of the CEFRLex project.

777,000 words distributed across several textual genres.

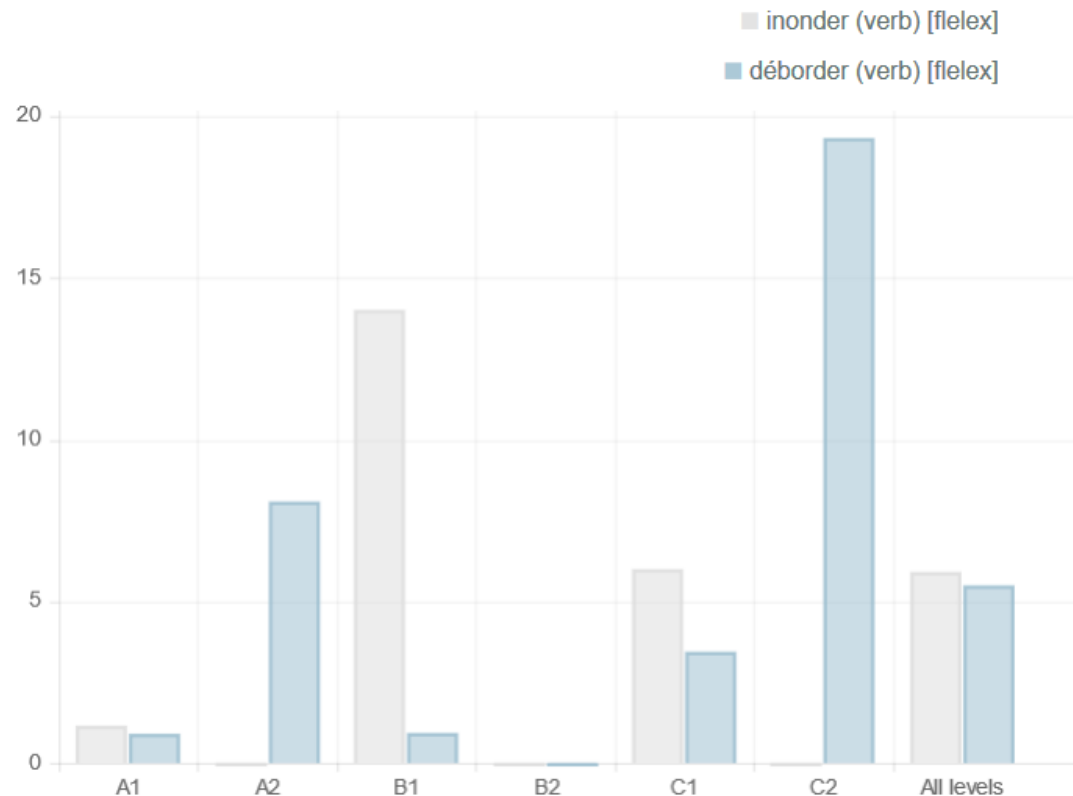Available online, possibility of comparison between 2 words.

Possibilities of using the FLELex resource: evaluating lexical simplification and vocabulary knowledge for learners of French

http://cental.uclouvain.be/flelex/

**Enter a word**

| inonder | déborder | – | Search |

Frequencies by CEFR levels for the words inonder* and déborder**.

- inonder (verb) [flelex]
- déborder (verb) [flelex]

# Analyzing language complexity

Gala, François, Bernhard & Fairon (2014)

How complex is a complex word ?
Complexity (objective) / difficulty (subjective)

Orthography
- Length (phonemes, letters, syllables)
- Orthographical neighbourhood
- Grapheme-phoneme coherence
- Syllable structure

Morphology
- Length (morphemes)
- Frequency of morphemes
- Size of the morphological family

Semantics
- Polysemy



Example in French (theft ... burglary... robbery) :

*vol – fuite – attaque – effraction – cambriolage – chapardage – acte de brigandage - maraudage*

# Analyzing language complexity

Gala, François, Bernhard & Fairon (2014)

How complex is a complex word ?
Complexity (objective) / difficulty (subjective)

Orthography
- Length (phonemes, letters, syllables)
- Orthographical neighbourhood
- Grapheme-phoneme coherence
- Syllable structure

Morphology
- Length (morphemes)
- Frequency of morphemes
- Size of the morphological family

Semantics
- Polysemy



Example in French (theft … burglary… robbery) :

*vol – fuite – attaque – effraction – cambriolage – chapardage – acte de brigandage - maraudage*

# ReSyf: a lexicon with graded synonyms

Gala, François & Fairon (2013), Billami, François & Gala (2018)

https://cental.uclouvain.be/resyf/



Interface: **D. Ricci & B. Delmée** -(2017-2018)
supervised by T. François (CENTAL) & N. Gala (LPL)

# Modelling individual difficulties

Larmuseau, Cornelis, Lancieri, Desmet & Depaepe (2020); Tack (2021)

Aims:
- gauging individual overall cognitive load to process (read, understand) a word
- accounting for individual differences between readers
Implicit / indirect measures:
- reading times, eye fixations, physiological data (brain signal, heart rate, skin temperature)

Explicit / direct measures:
- vocalization (read-aloud), verbalization (think-aloud), self-assessment

Building graded resources which

- include individual indirect measures for grading vocabulary
- propose texts / exercises according to personalized needs

# ALECTOR: a parallel corpus

Gala, Tack, Javourey-Drevet & François (2020)

79 original and simplified French texts for reading training online.
Lexical simplifications: Manulex and ReSyf.

Grades according to the difficulty of reading (reading times gathered in 6 schools, 970 children, 2017 to 2019) Javourey-Drevet (2021)

Interface: **S. Lâm** (2019) supervised by C. Ramisch (LIS) and N. Gala (LPL)

Z score: mean of different readers

http://corpusalector.huma-num.fr/

# Using graded resources in the classrooms... or in personalized trainings

For the teacher, in addition to other activities for vocabulary learning:

- Analysing texts before using them in the classrooms, identify complex words for a given grade (choosing or discarding a text)

During the class:

- Discussing about the knowledge of a word within a grade (whether the word is understood, a synonym can be proposed by the group, the word can be re-used in another context, etc.)

- Studying the morphology, the syntactic properties and the semantics of the word (POS category, cooccurrences, synonyms or thematic links –if possible browse through the semantic links)

In total autonomy:

- Working with texts adapted to the student profile, adaptive learning (Kerr, 2016)

# Conclusions and future work

New field with high potential for educational applications.

CEFRLex project is a pioneer in graded resources development.



Methodological challenges:

- Model more fine-grained gradings (e.g. for multiword expressions and collocations, for domain-specific texts)
- Include cognitive data in the gradings (to go beyond frequency distributions)
- Train personalized models
- Include graded resources in language learning platforms (track learner's activities and propose adapted contents)
- Extend to a different languages and varieties (e.g. oral)
- …

Interdisciplinarity:  linguistics, education, NLP, cognitive sciences…

# Graded resources of tomorrow



Smart Education
Adaptive Learning

iCALL

Hybrid devices (distance/face-to-face)
Empowerment, Autonomy
Interaction

Embodied cognition
Eye-tracking
Brain activity
Physiological data

Education Sciences

Cognitive Sciences

SLA

Language Sciences

Data Sciences

IA

Word Senses
Collocations
Multiword Expressions

Quantitative data
Deep Learning
Language Models
Big Data

NLP

Computational Readability
Text Simplification

© N. Gala

20

# Take home message

Graded resources are new resources with high potential for educational applications.

Beyond frequency distributions, there are important methodological challenges requiring interdisciplinary expertise.

nuria.gala@univ-amu.fr

**PAROLE ET LANGAGE**
UMR 7309 · CNRS · AMU

**ILCB**
Institute of
Language, Communication
and the Brain

Aix*Marseille
université
Socialement engagée

*Thank you*

# References

Billami, M. B., François, T. & Gala, N. (2018) **ReSyf** a lexical resource with graded synonyms. Proceedings of 27[th] International Conference on COmputational LINGuistics *COLING* Sta. Fe, US, 2570 - 2581.

Collins-Thompson, K. & Callan, J. (2005) Predicting reading difficulty with statistical language models. *Journal of the American Society for Information Science and Technology* 56(13), 1148-1463.

Deutsch, T., Jasbi, M., & Shieber, S. (2020). Linguistic features for readability assessment. *arXiv preprint arXiv:2006.00377*.

Dubay, W. (2004) The Principles of Readability. *Impact Information*.

Figueras, N. (2012) The impact of the CEFR. *ELT Journal* 66(4). Special issue October 2012. Cambridge University Press.

François, T. (2009) Combining a statistical language model with logistic regression to predict the lexical and syntactic difficulty of texts for FFL. In *Proceedings of the 12[th] Conference of the EACL: student research workshop*, 19 – 27.

François, T., Gala, N., Watrin, P. & Fairon, C. (2014) **FLELex**: a graded lexical resource for French foreign learners. *International conference on Language Resources and Evaluation (LREC 2014)*, May 2014, Reykjavik, Iceland.

Fulcher, G. (2010) The reification of the Common European Framework of Reference (CEFR) and effect-driven testing. *Advances in Research on Language Acquisition and Teaching*. Greek Association of Applied Linguistics. Thessaloniki, 15-26.

Gala, N., François, T., Bernhard, D. & Fairon, C. (2014) Un modèle pour prédire la complexité lexicale et graduer les mots. In proceedings of the conference *Traitement Automatique des Langues Naturelles TALN'2014*, Marseille, France, 91-102.

Gala, N., François, T. & Fairon, C. (2013) Towards a French lexicon with difficulty measures: NLP helping to bridge the gap between traditional dictionaries and specialized lexicons. *Proceedings from eLex – Electronic Lexicography*, Tallin Estonia, 132 – 151.

Gala, N., François, T., Javourey-Drevet, L. et Ziegler, J.-C. (2018) La simplification de textes, une aide à l'apprentissage de la lecture. Dans *Langue Française « Lire – écrire : Des savoirs scientifiques aux savoirs pratiques »*, 199 (3). Éds. Liliane Sprenger-Charolles et Alain Desrochers. Armand Colin, 123-131.

Gala, N., Tack, Javourey-Drevet, L., François, T. & Ziegler, J. C. (2020) **Alector**: A Parallel Corpus of Simplified French Texts with Alignments of Misreadings by Poor and Dyslexic Readers. *Language Resources and Evaluation for Language Technologies (LREC)*, Marseille, France.

Javourey-Drevet, L. (2021) La simplification de textes comme outil pour améliorer la fluidité et la compréhension de lecture chez les enfants à l'école primaire. Une étude en longitudinal avec des textes littéraires et scientifiques chez des enfants entre 7 et 9 ans. Thèse de doctorat, école doctorale Cognition, Langage, Education. Aix Marseille Université.

Kerr, P. (2016) Adaptive learning. Technology for the language teacher.

Larmuseau, C., Cornelis, J, Lancieri, J., Desmet, P. & Depaepe, F. (2020) Multimodal learning analytics to investigate cognitive load during online problem solving. *British Journal of Educational Technology* vol 51 Issue 5, 1548-1562.

Lété, B., Sprenger-Charolles, L. & Colé, P. (2004) **Manulex**: A grade-level lexical database from French elementary-school readers. *Behavior Research Methods, Instruments & Computers*, 36, 156-166.

Martinc, M., Pollak, S. & Robnik-Šikonja, M. (2021) Supervised and unsupervised neural approaches to text readability. *Computational Linguisticc Journal* volume 47 Issue 1 available online at https://doi.org/10.1162/coli_a_a00398

Tack, A. (2021) *Mark my words. On the automated prediction of lexical difficulty for foreign language readers*. PhD dissertation. KU Leuven & Université catholique de Louvain.

Tack, A., François, T. & Fairon, C. (2016) Evaluating lexical simplification and vocabulary knowledge for learners of French: possibilities of using the **FLELex** resource. In *Proceedings of the 10th International Conference on Language and Resources Evaluation (LREC 2016), Portorož, Slovenia: European Language Resources Association*, 230–236.